# The Digitization of Historic Newspapers on Microfilm: The Kentucky Experience

*by Kopana Terry*

*Kopana Terry (kopana.terry@ uky.edu) is Program Manager, University of Kentucky, National Digital Newspaper Program.*

As the Kentucky representative in the National Digital Newspaper Program (NDNP), the University of Kentucky Libraries Preservation and Digital Programs (UKPDP) team has worked extensively with historic newspaper digitization from microfilm over the last four years, using both an in-house production methodology and vendor resources. With more than 50 years experience with microfilming newspapers added to that, UKPDP is well versed with issues related to historic newspapers on microfilm. "Digitizing historic newspapers from microfilm" may sound as if all the work lies in the mechanics of digitization. Our experience tells us otherwise. If the digital surrogates are to be an accurate representation of the *newspaper,* there are several points to consider beforehand that have little or nothing to do with the digitization itself but, rather, with the newspapers and how they were microfilmed. This article identifies the more pressing of these issues and offers some solutions for them. It does not address in detail the more complicated affair of hardware, software, interface access, or storage associated with the digitization.

## A Brief History: Newspapers

The industrial revolution of the mid-nineteenth century ushered in the widespread use of acidic papers. Acidic paper is very volatile (i. e. brittle); unfortunately the majority of newspapers are printed on it, making the paper itself of little historic value and the content extremely imperiled.[1]

Given the temporary nature of acidic newsprint, it is not unusual to find newspapers that are little more than fragments. Trying to identify these frag- ments is akin to forensic science. In some cases, they have no text that confirm a date or title with 100 % certainty. Sometimes ads, ad placement, or serial articles are the only way to make a positive identification, except for an obscure date or notable reference that requires reading the text to discover. Pages closest to the outer covers of bound volumes tend to fall into this category.

Whatever the condition of a page, it is not the directive of librarians, archivists, or microfilmers to decide what may, or may not, be of value to a user. Microfilming the fragments was, and still is, a primary objective of content preservation in both microfilm and digital formats.

## A Brief History: Microfilm at the University of Kentucky

Without the mass of historic newspapers on microfilm that

we have today, America's history would largely be lost to the cobwebbed attics and wallpapered halls of American lore. An issue or two might have survived, but without early microfilming efforts, the most comprehensive record of our shared history would be as lost to us as if it had burned in the Library of Alexandria. Unlike precarious newsprint, polyester-based silver-halide microfilm is a proven preservation medium that can last up to 500 years under the right conditions.[2] There is no other media – save for rag paper, desert-bound papyrus or perhaps clay tablets – that can tout that kind staying power.

Microfilming of newspapers started in the 1930s at Harvard and Yale Universities, New York Public, and the Library of Congress.[3] Microfilming at the University of Kentucky began a decade later through the efforts of historian and UK Professor Dr. Thomas D. Clark and library director Dr. Lawrence Thompson. The two pioneers traveled the state – portable Recordak camera in tow – to microfilm holdings courtesy of newspaper publishers. Their lighting was horrible. The four-corners of the paper were so dark that the text was illegible while the center so bright the text was all but washed away. The focus was questionable at best. But however primitive their methods may have been, the two men microfilmed information before it was lost or destroyed.

By 1955 the UK Libraries had established an onsite newspaper microfilming operation. This *greatly* improved quality but, make no mistake, standardization was a long way off yet.
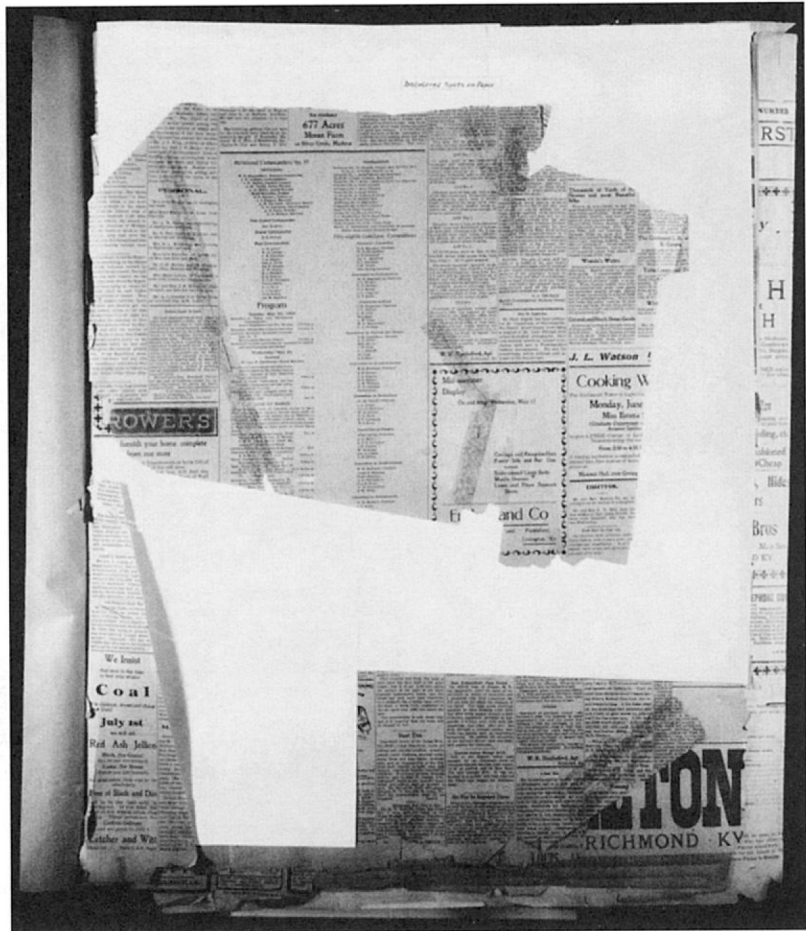


**Figure 1. Example of a page fragment**

Instead, home-grown policies were adopted; a practice we call "voodoo" microfilm.

It was during this period that microfilming took on individual characteristics. Some microfilmers methodically placed side date targets with each page, while others targeted only the first page of an issue. More amusing was the use of a glass ashtray – with burning cigarette – to hold down the corner of a page. Other microfilmers used pencils. They used whatever was available to get the job done. However, when the microfilm is digitized, objects that cover text can inhibit optical character recog-

nition (OCR) or they may cause detection errors in a scanner.

In 1981, UK Libraries became one of the first five institutions to participate in the United States Newspaper Project (USNP). During this time nearly 5,000 titles were cataloged and UK's microfilmed newspaper pages increased by an additional 1.5 million pages.[4] To date, the UK Libraries' vault houses nearly 30,000 reels of master negative film.[5] UKPDP continues to microfilm more than 150 current Kentucky newspapers as well as historic newspapers as they surface. After all, the *only* difference between "current"
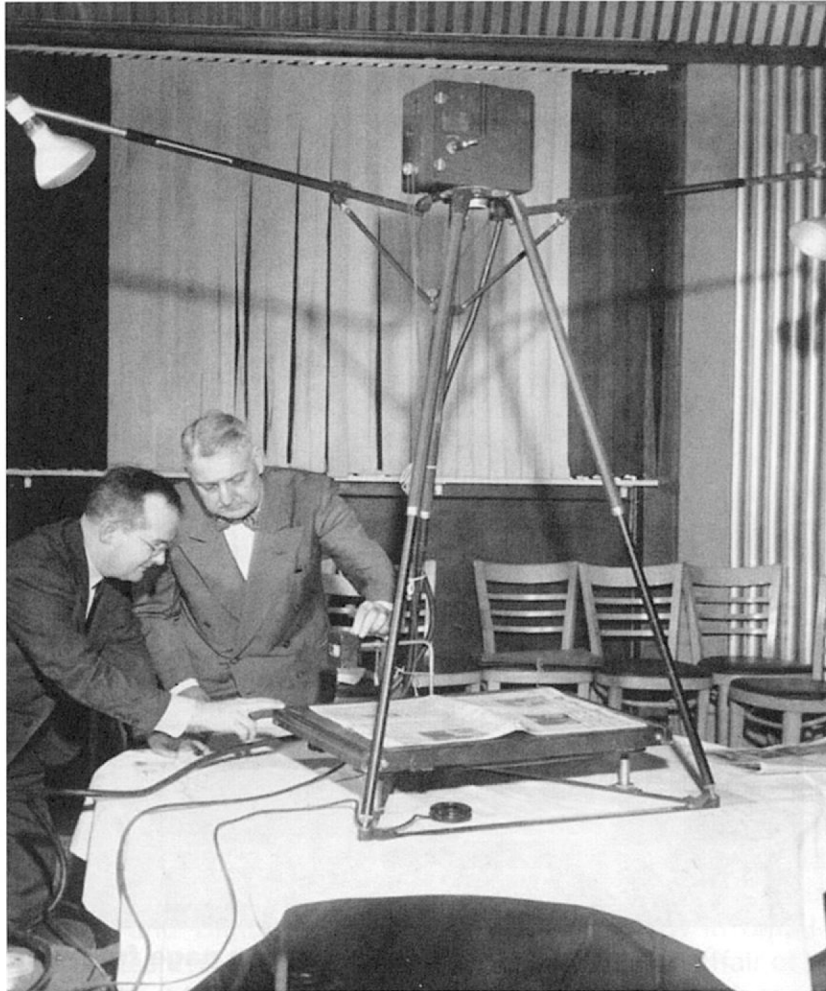
**Figure 2. Drs. Thomas D. Clark and Lawrence Thompson with their portable Recordak camera (courtesy University of Kentucky Archives)**

and "historic" newspapers is where they exist on a timeline.

To make this all possible, UKPDP has one of the last surviving full service microfilm labs in a library. This rarity has proven extraordinarily useful for newspaper digitization. We can choose a master negative that may have less than stellar density readings, and our darkroom managers can adjust duplicating procedures for remarkable improvement of the print master. This enables UKPDP to deliver better

microfilm and digital products.

Libraries without microfilm facilities were forced to use commercial vendors for these same preservation services. Ownership of that film has become a hot topic now that digitization from microfilm has come to fruition. Although many libraries may have retained ownership of the content, many micropublishers are asserting ownership to the film, if for no other reason than to recoup storage fees.

## History Repeating Itself

When microfilming operations began, the basic premise was to microfilm as much newspaper content as could be found. The rationale was to get the newspapers on film in any arrangement and let the user sort them out so long as the content was preserved on film. Often there was little forethought put to arrangement beyond very basic chronology. Standards and guidelines evolved as programs developed. USNP saw disparate programs around the nation come together under a single umbrella using a cohesive set of guidelines that were both meaningful and simple. Because some of the content we find compelling to digitize today was made during or before USNP standards took root, we sometimes encounter microfilm with little logic to the organization of the content.

In some ways we run the risk of history repeating itself with newspaper digitization. Those many reels of "let the user sort it out later" have come back. Later is now. We have the choice to either reproduce the mistakes that were made when the newspapers were microfilmed or remove the mistakes so that users of the digital content get as honest a surrogate of the newspaper – *not the microfilm* – as possible. The UKPDP microfilm-to-digital methodology is a concerted effort to keep history from repeating quite so literally.

## Microfilm Evaluation

We knew going into NDNP that to effectively digitize what was

on the film(s), we had to first know exactly *what* was on the film(s). Publishers make mistakes. Binders make mistakes. Microfilmers make mistakes. We all make mistakes!

When we choose titles for digitization, we evaluate each reel of that title from beginning to end. We believe that the microfilm evaluation step of our workflow is the foundation for the quality of all subsequent steps in the process. We store the evaluation information in a MySQL database that is continually accessed throughout the digitization process.

For a single two-year phase of NDNP we evaluate approximately 150 – 175 reels of microfilm. A time intensive process, microfilm evaluation is not easy if short staffed. It also is not particularly easy to train new employees or students to do high-level evaluation since it requires focused attention to detail.

There are three simple words we use when looking at newspapers on microfilm:

- Collation – to arrange in proper sequence/to verify arrangement

- Completeness – having all parts or elements

- Collection – the amount of material in one location

We think of this as the 'you have to know where you've been to know where you're going' method. If you take nothing else from this reading, take this: you cannot assume that everything you think should be on a microfilm reel really is. No matter what the date range or description says on the reel box, inventory, or database, you do



**Figure 3. This is an example of light-drop-off that is seen in some of Thompson and Clark's microfilm.**

not know what is on a reel until you have looked at it.

## Collation:

In our shop, during the process of microfilming a newspaper, we collate each issue to assure quality, proper order, and correct dates. We also inspect the master negative film for physical problems and bibliographic integrity. Mistakes are corrected (refilmed) so that the master negative is as bibliographically complete and physically pristine as we can make it. With older microfilm we don't have such luxury; the newspapers are no longer in our possession.

During collation of older microfilm we note incorrect or questionable dates, mispaginated, duplicate, and out-of-order pages, chronology, and any other peculiarities that can cause problems during digitization.

Such peculiarities might include a weekly newspaper that printed two issues on a single day. We see this most often during holidays or special events, like the assassination of Kentucky Governor William Goebel in 1900. It is not unusual when both front pages look strikingly similar but have no printed edition label. To recognize that these are not duplicate issues, we inspect article headlines and then, hopefully, find some clue as to which issue came first.
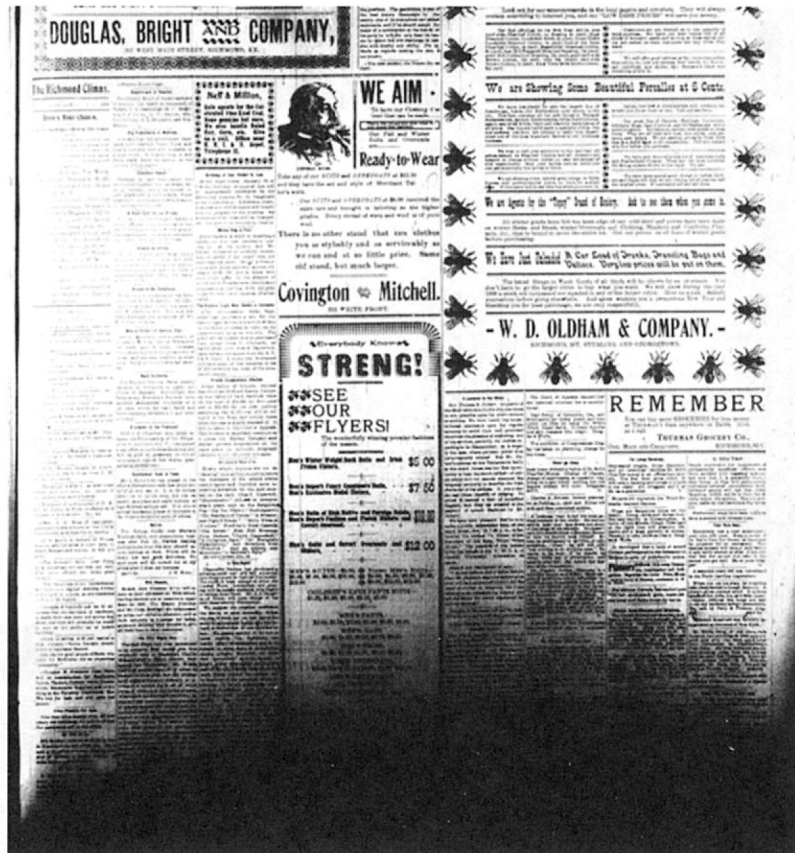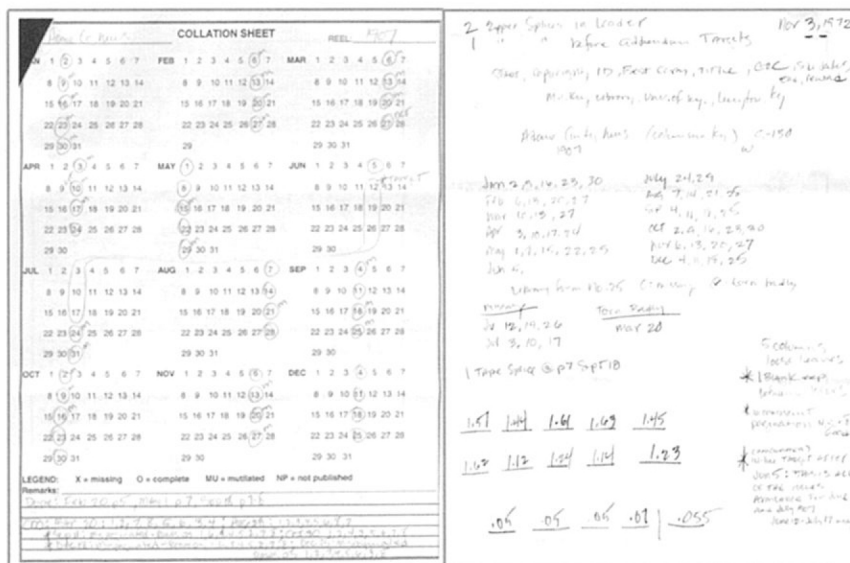
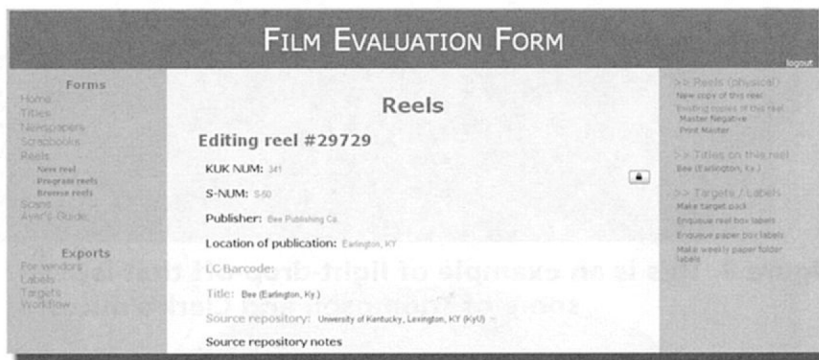Collation by hand (above) transcribed into a MySQL database (below)

**Figure 4. Collation sheets with MySQL interface**

During digitization the issues will be intertwined unless edition labels are assigned. For instance, without edition labels the user will see pages 1,1,2,2,3,3,4,4 instead of seeing *Morning Edition* pages 1,2,3,4 before *Evening Edition* 1,2,3,4. This lack of organization could adversely impact research.

Another advantage to a thorough collation is finding publication patterns. When looking at a calendar, noting days with present content makes missing issues extremely easy to spot. Inspecting the issue numbers can further identify issues that

were missed during filming or simply not published. The latter is often the case on or after the Christmas and New Year holidays (see Figure 4).

To be sure, digitization software can fix most sequential errors; sorting by year, month, day, and page number as desired. Software may also order by reel sequence number – a provenance marker assigned as part of the NDNP specification. But as we have established, not everything is in order on the film, so sorting by sequence number has the possibility to repeat those mistakes.

Noting chronological order can also signal something bigger. We have found issues from 1893 on film with a start date of 1894. These issues were missing from the 1893 reel. So, we actually have issues that we previously thought were missing. This kind of one-off issue on a reel is not bad or even wrong. Quite often, lone issues will trickle in and then are microfilmed with other issues from that title. Or, over time, stray issues from multiple titles are compiled and filmed as, what we call, "miscellaneous" reels. The linear nature of analog materials can create confusion for the user. But a digital interface can bridge linear restrictions such that, no matter where an issue came from, the corpus of a title can be coherently viewed.

It should be pointed out that, from a digitization vendor's standpoint, it is not their prerogative to make the kinds of decisions that are discussed here. They are paid to digitize what is on the film with little intellectual judgment beyond repeating obvious dates, page numbers, and edition labels. A benefit to doing this kind of work in-house, is that we have found that our evaluation of the microfilm is invaluable when we work with digitization vendors. We supply guidelines that explain what is on the film and how their metadata technicians should handle noted special circumstances. To our knowledge, no other digitization institution supplies such detailed guidance.

**Completeness:**

During the microfilm evaluation process we inspect each reel for

missing pages, missing issues, and unpublished issue dates. Some of these errors are due to the microfilmer. For instance, they may have turned two pages rather than one. Some are due to binding errors – volumes were usually bound exactly as the publisher presented them, with duplicate and out of order issues and pages. Some are due to publishing errors that may be a hundred years old. It is not rare to find the wrong year printed on a first issue in January. The digitization software is very sophisticated and can correct for most of the problems, but it cannot replace missing pages and issues, and it cannot decide what is or is not a correct date.

A unique problem often overlooked is that of unpaginated and mis-paginated pages. Advertisers at the turn of the century had a tendency to use the same ads over and again. Publishers would also print them in the same location on the same page from issue to issue. The same can be said for serial articles like agriculture updates and social happenings. This predictable printing tactic turns out to be advantageous when evaluating unpaginated and mis-paginated pages. Once a pattern has been established using these repeating ads and articles, out of order unpaginated or mis-paginated pages are easy to spot.

Of course, the simplest way to recognize page-two from page-three, for instance, is to examine the edge of the paper. If a rip, tear, mutilation, or even binding holes mirror that of page-one, it is very likely page-two. If these two ID tactics do not work,
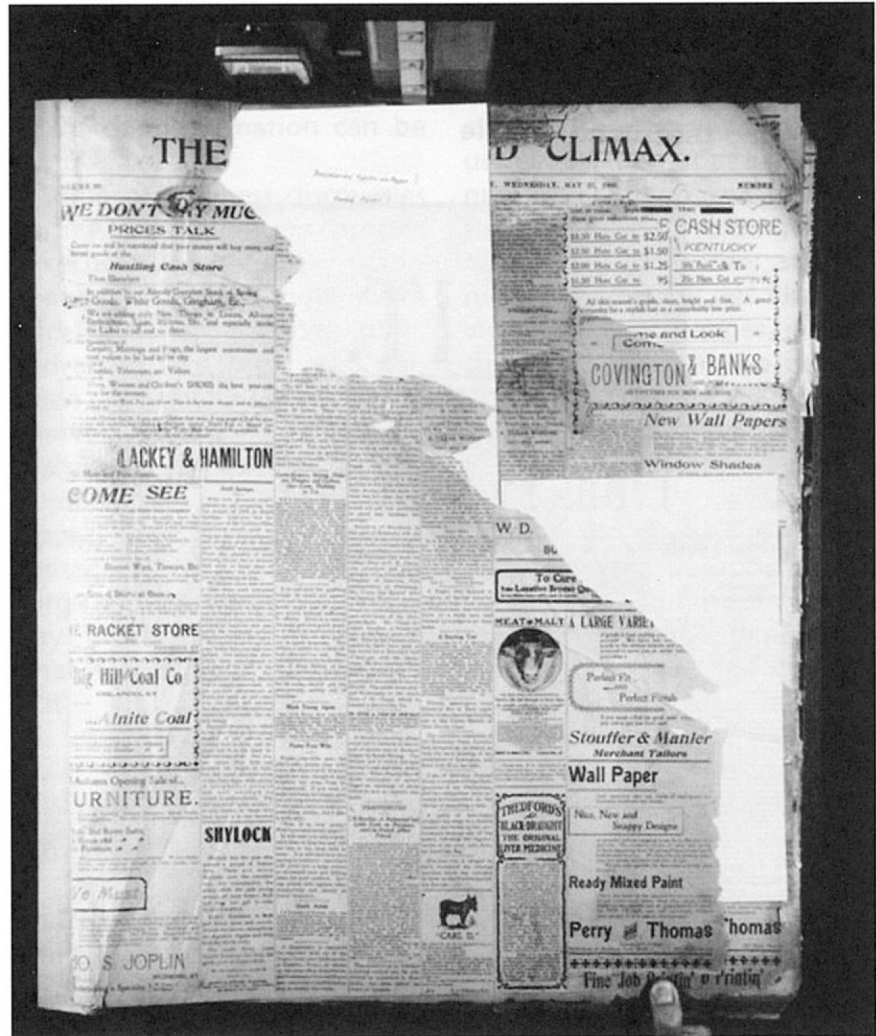


**Figure 5. A compilation image of a single page compiled from eight copies.**

the last thing left to do is to read the paper. Even then there is no guarantee that it can be perfectly identified. Sometimes you just have to guess.

A final word about page fragments. We do our best to logically identify fragments because, unlike microfilm, which is linear, digital files can be sorted (viewed) in any number of ways as mentioned earlier. In the case of fragmented pages, assigning the correct date is critical from a user's perspective. It would be confusing to find a page frag-

ment from May 8, 1893 at the beginning of the April 28, 1900 issue. That presentation would not be an accurate representation of the newspaper, which is why we put such emphasis on identification.

Vendors do not determine bibliographic integrity if it is not printed on the page. A very common habit we find is an issue with a four-page supplement in the center of its regular four-page issue. By today's standards that supplement would be positioned at the end
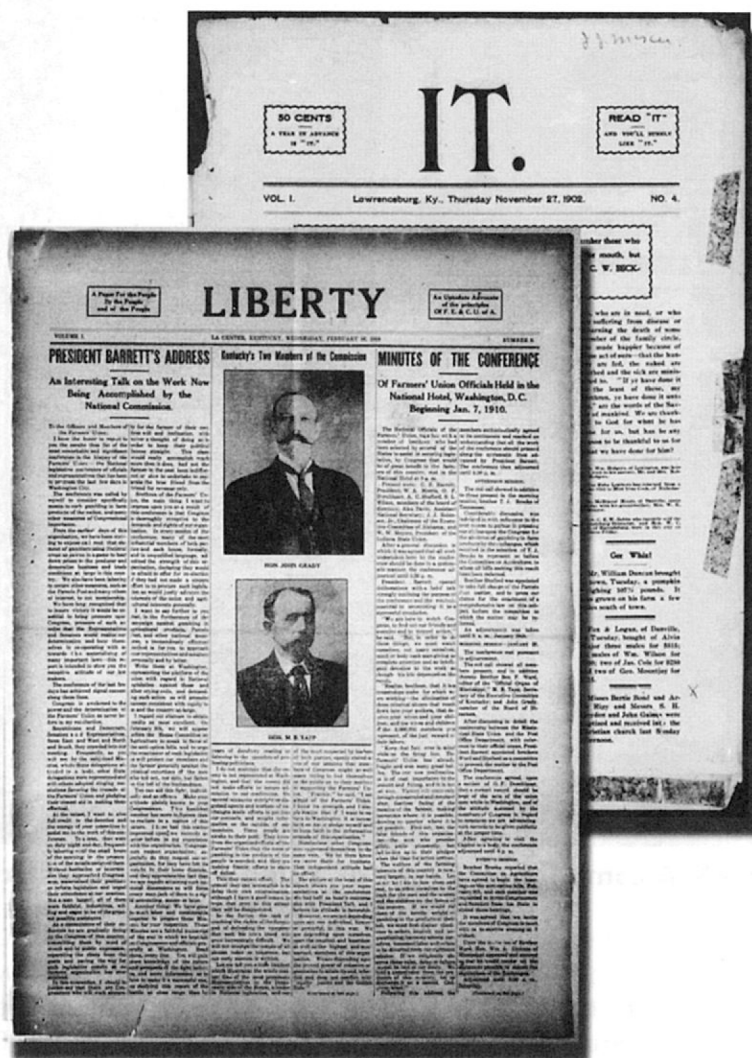
**Figure 6. A masthead sample**



**Figure 7. Examples of orphan papers**

seen microfilmers ignore publisher printed pagination and simply written in their own.

**Collection:**

There are many things to collect during the microfilm evaluation process. How much to collect is up to you but be warned; it is easy to get bogged down in too much detail.

We collect general information about mutilated pages/issues. Since most historic newspapers suffer some amount of broken page edges, we only consider pages truly mutilated if they are missing a great amount of text. The page fragments discussed previously fall into this category.

One of the most unique instances of mutilation that we have found to date is a single issue from 1908. The publisher kept, and the microfilmer photographed, eight copies of the first two pages (i. e. the first sheet, front and back). Each copy was missing text in different places. That meant each page housed varying degrees of the original information. When the first image of page-one was discovered during the evaluation process, it appeared to be nearly whole. Closer inspection revealed portions of all eight copies. The situation begged the question: do we digitize and present this compilation image and disregard the other seven – the microfilmer did not provide a compilation image of page 2, incidentally – or do we simply present all 16 pages and let the user deal with it? In this case, we opted for the latter because each page really did present the user with different content. To remove any one page would

of the issue. Without pagination, mis-pagination (often a result of well meaning microfilmers writing in page numbers), or a "supplement" label, the user will see the issue as 1,2, supplement pages 1 – 4, 3,4. If the microfilmer wrote the page numbers sequentially, then both the newspaper pages and the

digital files will appear as pages 1 – 8. The problem may be that content continued between pages 2 and 3, for example, will now be separated by four additional pages. A user turning from page 2 expecting to find the story on page 3 will have to scour the bulk of the issue to actually find it. Worse cases have

have been to deny a user of that content.

There are other, far less eccentric duplicative cases to be understood. For NDNP we are charged to digitize everything on a reel, including the duplicative material. Each page image, including all targets and empty camera bed exposures, become part of the reel sequence numbering system. Therefore, each exposure must be accounted for in order for that location/sequence number to be accurate. In theory, this will enable quick rescanning if anything should go wrong with the original data. Once the sequence numbering is established, one can then choose the duplicative material to keep, then discard the rest. This would lessen the burden of tape or server storage and possibly lessen any confusion it might cause a user.

In addition to the technical information we have discussed, intellectual information about a newspaper can be collected during the evaluation process as well. NDNP requires that we compose a 500 word historical essay for each title. Some of what we know about a newspaper we find through research of historical writings, but a good deal of what we learn about the newspapers we gather from the papers themselves. This could be anything from an editor's political leanings to the newspaper-printing house that burned and stopped publishing for two weeks. Some historic papers have remarkably ornate mastheads along with slogans that are quite provocative: *All The News That's Fit To Print; Official Organ of The Party In The Fourth Congressional District;*

*By Industry We Thrive; A Weekly Journal Identical In Interest With Its Own People,* and so on. All of this information can be useful.

Some of the best discoveries during the evaluation process have been the discovery of orphaned titles. Though every effort was made during USNP and again just three years ago when every reel in our vault was re-inventoried, re-boxed, and re-shelved, we still occasionally stumble upon orphan titles that slipped under the radar. The La Center *Liberty* and Lawrenceburg's *It* were both found this way. You don't know exactly what is on a reel until you look at it.

## Other Helpful Things to Know

### Multiple Title Changes for one Newspaper

One of the most challenging dilemmas is that of title changes within a single microfilm reel. Different from a miscellaneous reel with different titles from different physical places, some reels include variants of the same parent title such as *The Paducah Sun, The Paducah Weekly Sun,* or the *Paducah Sun Weekly Edition,* which are all variations of the same title. To the publisher, binder, and microfilmer, it was all the same newspaper so they made no effort to separate them by title. With the *The Paducah Sun,* for example, five title changes occurred in the span of a single decade. These bibliographic anomalies also appeared on multiple reels of microfilm. The multiple changes

were not discovered until after the title had been chosen for digitization, the film duplicated, and the evaluation process underway. It is not at all uncommon to have two or even three title changes on a single reel of film. The inventory and the microfilm box will not likely divulge the title changes, leaving the discovery until the evaluation process.

However, until Phase 2 of NDNP, we were charged to write an essay for every title change. These variations have obvious implications for the historical essay, and they forced an adaptation in the digitization workflow. When title changes occur within a NDNP reel, all digital end products, including metadata, must be delivered separately to the Library of Congress.

### Physical Characteristics and Microfilm Scanners

Microfilm scanners are designed to automatically detect a page edge and scan accordingly. We know this to be a flawless operation with our current-made film: The pages are evenly spaced and placed squarely on the camera bed with ample border on all four sides. They have even lighting on a high contrast camera bed so the page is strikingly visible with even densities from beginning to end. For instance, today's 100' reel can be scanned in under thirty minutes and produce approximately 600 uncompressed 400 DPI page images.

For historic newspapers on older film, sometimes few of these principles apply. The newspapers may be badly deteriorated

or poorly filmed such that even the best detection software will skip over the page (we're back to those page fragments again). They are rarely spaced evenly on the camera bed or even between exposures (controlled by the camera's "gate"). The reduction ratio may be so low that the page edge is just inside the exposure's edge or, perhaps worse, the reduction ratio may change throughout the reel without warning. The camera bed may be very close to the same color as the paper, making detection of the page edge by a scanner very difficult. If the film was over-exposed, detection could be nightmarish if not impossible, thus forcing the technician to scan one page at a time, a painfully slow process. The same is true of under-exposed film on a black background, though not quite to the extent of over-exposure.

This brings us back to collating the film. If all one sees are the digital files, one has no way of knowing if the scanner skipped a page or if it was never there. We have found that it is more time efficient to know what is on the film before it is scanned. This allows the scanning technician to recognize that a page has been skipped, stop the scanner, roll back the film, and scan the skipped pages. To do so after a reel has been scanned can mean re-spooling the scanner, fast forward or backward to the exact page – with some risk to the film just by touching it again, then following the correct numbering to insert new images into the group. It opens up a new avenue for human error that can be avoided entirely.

## OCR Accuracy

Readability of fragmented pages, bound volumes, and skewed pages can all present challenges to optical character recognition (OCR) accuracy. Nothing can be done about the first two problems. Obviously, it would be dishonest to fill in missing data of a fragmented article or text hidden by the tight gutters of a bound volume. In fact, early digitization efforts did precisely that, calling it "boutique" digitization. The practice is now widely regarded as unfruitful and not in the least bit productive. It improved OCR accuracy, no doubt, but there is no way of knowing if it increased search results for users.

In cases where a large bound volume has created the distinctive hump near the gutter such that the text is no longer straight on the page – not to mention the 'hot spot' caused from the intensified lighting – software that can correct for it is not typically used in newspaper digitization workflows. The mechanism used to flatten the page image is generally considered a 'manipulation' of the image itself and, therefore, an unacceptable practice.

To be clear, the simple act of digitization is a manipulation of the original object. The same is true of microfilm. Nicely detailed photographs, for example, will lose some amount of detail because of the high-contrast nature of the microfilm. Likewise, there are as many densities, contrast, and lighting possibilities in digital scanners as there are eyes to see them. To say that we don't manipulate the digital images is not entirely

truthful. Manipulation and change is the inherent nature of reformatting no matter the medium. The point, however, is to introduce as little manipulation as possible into the master image file, whether that be analog or digital. The master image, if left unmolested, can survive evolving technologies with as much of the original information intact as possible.

For now, if better OCR is a top priority, a work around of these master file restrictions is to make a copy or "work" file of each newspaper page image. Any number of manipulations can be applied to that work file, such as increased contrast to enhance the text, sharpening which will also enhance the text, or page leveling as described above. These techniques will likely produce better OCR. Most newspaper digitization institutions do not take this copy file route because of the sheer workload it would add to the already cumbersome workflow. Plus, there is still no way of knowing if it improves OCR accuracy to the degree that there is an increase of search result accuracy.

All that said, deskewing the newspaper pages is one "manipulation" that is allowed by NDNP. Deskewing is absolutely a must for OCR accuracy; not just for the text but to deliver correct column read order of a newspaper. Read order has been a significant cause for delay in newspaper digitization. 'Zones' that outline each column must be created to tell the OCR software where one column ends and another begins. Without zones, the OCR software will read left to right as if the page

were from a book, combining incongruent sentences from the different columns into a singular string. Words hyphenated from one line of text to the next would not be joined but be left as two distinctly different parts.

Zoning is an especially important feature when generating article level data rather than page level data as we do with NDNP. Articles often carry over from one page to the next or from one column to another. Advanced zoning connects the disparate parts of one article into a single image. It would be disastrous to the user to be offered conflicting stories in a single "article" image.

At any rate, unlike fragments and bound volumes, skew can be successfully corrected during the scanning process – most professional microfilm scanners can deskew the page images while scanning. Deskewing newspaper pages can also be performed after scanning with a deskewing application. Some of these after-scanning applications can be automated.

## Orientation, Reduction Ratio, and True DPI

Differing orientations and reduction ratios on microfilm can be particularly troublesome during scanning. Neither is unusual on miscellaneous reels where single issues have been cobbled together. But make no mistake, the same has been found within single title reels as well. Like so many other issues we have identified for consideration prior to digitization, the orientation (position) of the newspaper page on the microfilm should be added to that list. It takes little effort for users to turn their

reader this way and that in order to read a piece of microfilm. The film readers are designed to accommodate A and B and 1up and 2up positions alike. Scanners do not comply quite so easily. Plus, orientation will supply an approximation of page images on any one reel of microfilm. Therefore, it is important to know to fulfill a page count.

A major consideration during the scanning process is true DPI. Scanner manufacturers will tell you their products can scan "up to 400 DPI" or "up to 600 DPI". From their standpoint, that may be true. (They do not necessarily consider interpolation a problem, though we do and so do NDNP specifications.) Simply set the scanner software to capture at 400 DPI, make a grayscale scan, open the image in any image editing software and it will tell you "400 DPI". The devil is in the details, however.

True DPI is measured by the width of a page in pixels divided by the physical dimension of the width. A physical page that measures 15x20 inches and then scanned measures 4512 pixels in width is equal to 300.8 true DPI. (4512px/15"=300.8 DPI) If a scanner is set to capture at 400 DPI, this digital page image falls far short.

If a newspaper was filmed at a reduction ratio that is too high, above 20x for instance, it can prevent the digital page file from reaching 300 DPI, which is the lowest DPI acceptable for NDNP. (Incidentally, 300 DPI has, in the last few years, been quietly adopted as the minimum DPI for most digitized objects, not just newspapers.) Using an orientation that is unsuitable for the size of the original news-

paper coupled with a low reduction ratio will produce similar results. Either the paper will be too big on the film and cut-off, or nearly so, or the reduction is too high and a minimum DPI cannot be achieved. The point is to hope that the microfilmer used the correct orientation and reduction ratio for the newspaper in order to achieve true DPI.

## Closing

This article highlights only some of the unique aspects that UKPDP has found while digitizing historic newspapers on microfilm. Using the Three C method – *collation, completeness, and collection* – during microfilm evaluation, most filming and binding errors and other unusual anomalies can be identified. Knowing what those shortcomings are can keep mistakes of the past from being perpetuated into the future and yields predictable results with the digital end products.

## Endnotes

[1] Lisa L. Fox, ed., *Preservation Microfilming: A Guide for Librarians and Archivists,* 2nd ed. (Chicago: American Library Association, 1996).
[2] American National Standards for Imaging Media – Processed Safety Photographic Films – Storage, ANSI/NAPM IT9.11 – 1993, 5.
[3] Fox, Preservation Microfilming.
[4] "National Digital Newspaper Program: The Kentucky Edition," http://www.uky.edu/Libraries/ndnp/kyhistory.html.
[5] "USNP Preservation Microfilming Guidelines," http://www.loc.gov/preserv/usnpguidelines.html.