

# Bringing the Past to the Present

by John Herbert and Karen Estlund

**N**ewspaper digitization is exploding. As we move through the first decade of the 21st century, the number of digital newspaper initiatives is growing rapidly, which brings the past firmly into our electronic present. This phenomenon isn't limited to the U.S.—the International Coalition on Newspapers (ICON) cites more than 40 countries with digital newspaper projects (<http://icon.crl.edu/digitization.htm>). In the U.S., the National Digital Newspaper Program (NDNP), launched by the National Endowment for the Humanities (NEH) with assistance from the Library of Congress (LC), has created great interest in digital newspapers among academic and public institutions, as well as numerous commercial entities ([www.neh.gov/grants/guidelines/ndnp.html](http://www.neh.gov/grants/guidelines/ndnp.html)).

The J. Willard Marriott Library at The University of Utah, partnering with Brigham Young University (BYU) and Utah State University (USU), runs a pioneering statewide program that has changed the face of newspaper research. Since its inception in 2002, the Utah Digital Newspapers (UDN) program quickly became a model for other institutions across the country. We have been and remain on the leading edge for newspaper digitization, especially within the public sector.

Our program presents digital images and searchable text of historic Utah newspapers to the general public through our website (<http://digitalnewspapers.org>). Anyone with an interest in history can search by keyword or browse by title and date, from any computer. UDN is easily accessible, easy to use, and free of charge, and it is beginning to make reading Utah newspaper microfilm obsolete.

## HISTORY OF THE UTAH PROGRAM

In 2001, the Marriott Library received a Library Services and Technology Act (LSTA) grant to digitize 30 years of three weekly newspapers. With that funding, we launched the UDN website in December 2002, with 10,000 pages each from the *Wasatch Wave* (Heber City), the *Times Independent*

(Moab), and the *Vernal Express* (Vernal). This initial implementation proved an immediate success.

In 2003, the library received a second LSTA grant to continue digitizing newspapers. The library hired a project director and added 106,000 pages. Later in 2003, in a watershed event, the Institute of Museum and Library Services (IMLS) awarded the library a 2-year, \$1 million grant. With this grant, we demonstrated that newspapers could be digitized and served over the internet to a wide audience on a large scale. We added more than 268,000 pages, distributed content to BYU and USU, and expanded coverage to include 27 of Utah's 29 counties.

Several Utah public libraries, assisted by UDN, donated funds and/or obtained other LSTA grants to digitize their local newspapers. These are just a few of those libraries: In 2003 the Weber County Library provided \$50,000 to digitize the *Ogden Standard* and its predecessors; between 2004 and 2006, the Park City Library and Historical Society raised \$15,000 for the *Park Record*; and in 2007, Myton City with the Duchesne County Library raised \$28,000 to digitize nine titles from Duchesne County. The Davis County Library established a relationship with the *Davis County Clipper* to fund the *Clipper's* inclusion in UDN. The library continues to fund its digitization narrowing the gap to the present day issues.

In total, we have raised \$2.5 million over the past 5 years. We have been and will remain an entirely soft-money program.

## NATIONAL DIGITAL NEWSPAPER PROGRAM

In 2005, the University of Utah was one of six institutions across the country to receive a 2-year grant for NDNP Phase 1. The program aspires to eventually include content from all states and territories hosted at the LC's Chronicling America page ([www.loc.gov/chroniclingamerica](http://www.loc.gov/chroniclingamerica)). The scope of the project is English-language papers, primarily scanned from microfilm, dating from 1836 to 1922. For states that participated in the NEH-sponsored United States Newspaper Program (USNP), which cataloged and

preserved newspapers on microfilm, this is a natural transition. States that did a majority of their microfilming during the USNP process generally benefit from more recent and higher quality microfilm images. Microfilm remains the preservation medium for NDNP, although we provide "preservation" quality—high-resolution, digital images—to NDNP as well. The catalog records created during USNP are available on the Chronicling America site, which is a great resource to track down where newspaper titles are available, whether in print, microfilm, or digital form. The NDNP content we generate will be available on the UDN website, but Utah content digitized from other projects will not be available from Chronicling America.

In Phase 1 we added 105,000 pages dating from 1900 to 1910. Today we are one of eight awardees in Phase 2 of the NDNP, in which we will process another 100,000 pages dating from 1880 to 1910.

### CURRENT HOLDINGS

At this writing, UDN houses 602,000 newspaper pages containing more than 6 million articles. It holds 50 titles, including the first issues of the *Deseret News* (1850) and the *Salt Lake Tribune* (1871), which today are the two leading newspapers in the state. The time period covered in the various papers is 1850–1969.

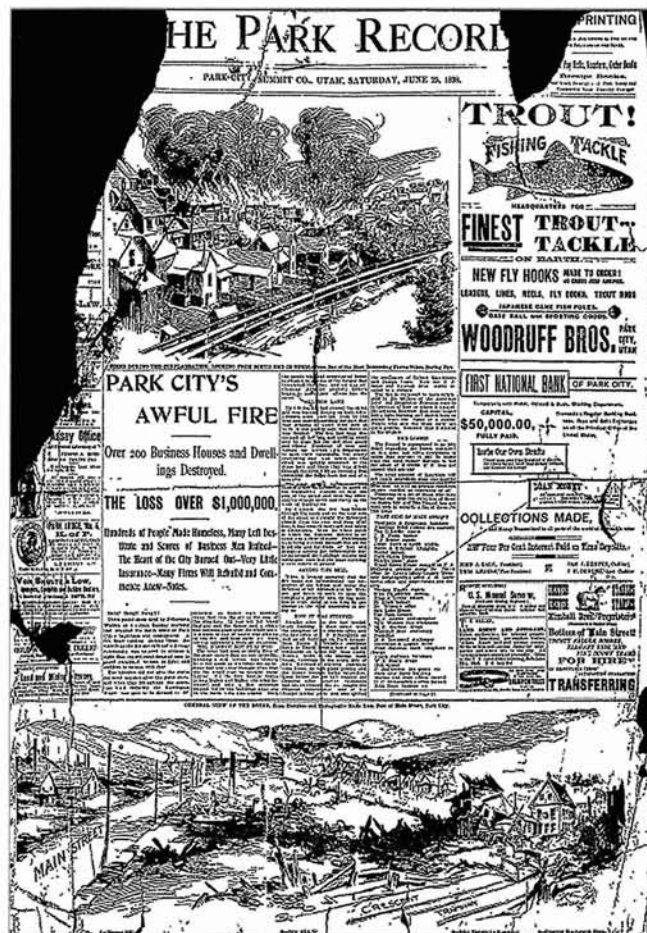
Website usage increased fourteenfold—from 60 visits per day in 2003 to 830 in 2006, when the program grew for 11 consecutive calendar quarters. It is amazingly popular with those who use it; 88% of readers rate us as "good" or "excellent." As you might expect with online content, 42% of our users are from outside Utah. A sampling of spectacular issues includes the following:

- *Park Record*, describing "Park City's Awful Fire" on June 25, 1898
- The many exploits of Butch Cassidy, one of which is related in the *Ogden Standard*, July 5, 1902
- *Eastern Utah Advocate*, Nov. 26, 1915, reporting on the death and burial of Joe Hill
- *Manti Messenger*, April 19, 1912, misreporting that on the Titanic "Three Thousand Passengers Have Narrow Escape From Watery Grave"
- *Vernal Express*, Dec. 11, 1941, declaring "Japan Attacks US" on page 3 (this just happened to be the paper's extravagant Christmas edition)

You can find the *Broad Ax*, an African-American newspaper run by Julius Taylor around the turn of the 20th century. Taylor later moved the *Broad Ax* to Chicago, where he earned some national acclaim. Also included is the *Topaz Times*, the newspaper from the World War II Japanese internment camp near Delta, Utah.

### USING THE UDN WEBSITE

Accessing content is as simple as connecting to our website. The entire collection is full-text searchable on the



Park Record description of "Park City's Awful Fire" on June 25, 1898

homepage. Individual titles can be searched or browsed from their own webpages. The user can browse by publication year, view all births, marriages, and deaths announcements, and perform specialized searches. One of the best and most popular ways to find regional news is with our county map. The user can view a map of Utah and select available titles in each county.

In April 2008, keyword searches found 10,170 results for **suffrage**, 8,422 results for **statehood**, and 161 results for **Butch Cassidy**. Users can perform exact phrase searches and limit to different fields or a particular group of titles. Searching for prohibition and either **women** or **ladies** in Davis County papers produces 202 results.

Readers increasingly use UDN as a source for various types of research. Val Holley, who lives in Washington, D.C., recently published an article titled, "Leo Haefeli: Utah's Chameleon Journalist," in the spring 2007 issue of the *Utah Historical Quarterly*. He listed 40 references to articles in UDN. Roy Webb, a staff member at the Marriott Library, presented the history of Utah soccer at the 2007 Utah Historical Society annual meeting. He based his research in UDN.

Genealogists make up our biggest group of readers. One of them, Thomas Billis, from West Jordan, found an article

## BUTCH CASSIDY IN A TRAIN ROBBERY

Utah Outlaw Believed to be a Member of the  
Kid Curry Gang That Held Up the  
Rock Island Express.

Chicago, July 5.—That members of the Kid Curry band of bank and train robbers, wanted for alleged complicity in the recent Union Pacific hold-up, perpetrated the robbery of the Rock Island express train at Dupont Thursday night, it is believed probably by detectives, one hundred of whom are working on the case today.

Charles Nessler, the boy who was stealing a rifle on the train when it was

stopped, described the men to detectives today and his description is said to tally with photographs and descriptions of Butch Cassidy and Sundance Kid, alias Harry Longbaugh, alleged members of the Kid Curry gang. It was officially stated by an officer of the United States express company, that the robbers secured only \$50 worth of jewelry. They carried away a package of worthless vouchers and other papers but overlooked a package containing \$100,000.

*The many exploits of Butch Cassidy, one of which is related in the Ogden Standard, July 5, 1902*

about himself from 1930, when he was 5 years old and got lost in a canyon near Price, "seeking body of lion." An Australian academic researcher used UDN for researching the early cement industry. Community theater guilds read advertisements to see popular fashions from the time periods their plays are set in. Even commercial firms use our collection—a food industry representative used UDN to research the history of tofu and soy food production in the West. The interesting uses of UDN grow every day.

### CONTENT SELECTION

With an estimated 8 million to 10 million pages of newspaper content in the state to digitize, content selection is paramount. The more practical considerations are availability, quality, and format of the source materials. UDN has an advisory board composed of librarians, historians, journalists, and publishers that help guide content selection. Some considerations include the following:

- Rural or metro—To be truly statewide, the project should extend widely to include content from all areas of the state. We have had great success generating broad interest by processing titles from rural areas first. In fact, there is so much interest across the entire state that many local areas are now raising their own funds to extend the collections we did initially with grant funds.
- Daily or weekly—Daily papers provide more news coverage and generally have a broader audience, while weeklies allow you to process a longer run. To put it another way, dailies are much more expensive to process because they have so many more pages. In our earlier stages, we preferred to focus on weeklies because, with a given amount of money, we could cover a longer time period and cover a larger geographic area of the state.

## HILLSTROM IS BURIED; FUNERAL TYPICAL ONE

CHICAGO, Nov. 21.—They're burying Joe Hillstrom tomorrow. It's a big event in Chicago's "rebel" ranks.

The ghetto, the slums, the lodging house quarters and the manufacturing districts are buzzing with preparations for the funeral of the Industrial Workers of the World song writer who was executed in Utah after conviction of murder. Anarchists, nihilists, saboteurs and secret organizations whose existence is based on opposition to present social and economic systems, are taking part in the honorary rites over the body of their dead comrade.

Handbills in a dozen different languages have been circulated, urging large attendance at the West Side Auditorium at 10 o'clock tomorrow morning. Today hundreds of persons passed into the undertaking rooms and viewed the body of the poet whose full name was Hillstrom.

An Industrial Workers of the World band of forty pieces will furnish funeral procession music. The Russian mandolin club will provide music during the services. Hillstrom's own songs will be sung and played.

Black bordered handbills tied with cord in the organization color of red and black have been circulated with the program for the exercises. In black type on the front of one is the inscription:

"In memory of Joe Hill, murdered by the authorities of the state of Utah, November 19, 1913. He died that men might live. We never forget."

O. N. Hilton, attorney of Denver, who defended Hillstrom and who represented the defense in the Haywood, Pettibone and Moyer trials, will deliver the funeral oration, which will be followed by addresses by James Larkin of Dublin, the man who led the dock workers' strike in England, and William D. Haywood.

*Eastern Utah Advocate, Nov. 26, 1915, reporting on the death and burial of Joe Hill*

- Length of run—The trade-off here is whether to have more titles with shorter runs or fewer titles with longer runs, and each approach has merit. In Utah we have tried to strike a balance between the two by digitizing approximately 10,000 pages the first time we select a title. For an eight-page weekly, this provides a 24-year run, while for a 16-page daily, it provides only 2 years.
- Editorial policy—Historic newspapers, like papers today, present varied editorial views. Our advisory board helps identify the editorial policies of new titles, supporting opposing views to help create a balanced collection.

### ORIGINALS VERSUS MICROFILM

When selecting content, another dimension to consider is whether the scanning source material will be original newspaper or microfilm. The principal reason we prefer using originals is that they produce new, high-quality digital images. This eventually leads to better online viewing and more accurate keyword searching. (For a more detailed discussion of paper versus film and the related optical character recognition (OCR) issues, see Kenning Arlitsch and John Herbert, "Microfilm, Paper, and OCR: Issues in Newspaper Digitization," *Microform and Imaging Review*, Vol. 33, No. 2, spring 2004, pp. 59-67.) Approximately half our volume is from originals.

One prerequisite for using originals is to find them in good enough condition to be transported and scanned. In



The Broad Ax, an African-American newspaper run by Julius Taylor around the turn of the 20th century

some cases we have utilized the services of our colleagues in the preservation department at the Marriott Library, University of Utah. They have done minor repairs such as surface cleaning and mending small tears. Some originals are simply in too poor of a condition to be easily repaired, and in these cases, we have opted to scan from film.

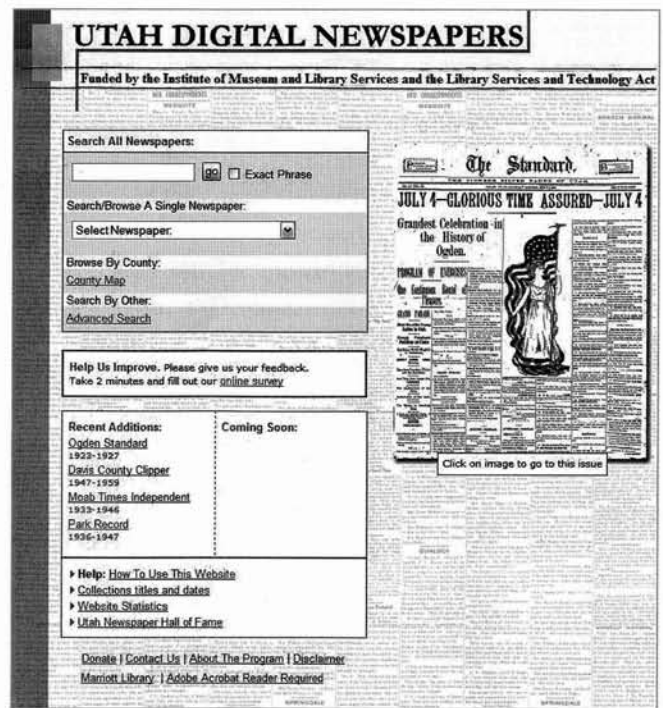
Using originals requires a local or in-state scanning service, since bound newspaper volumes are difficult to transport. It also requires an overhead scanning process, since large and sometimes fragile bound volumes cannot be scanned on a flatbed.

Alternatively, there are many advantages to using microfilm as the source material for scanning. In almost all cases, because of the success of the USNP, major newspaper titles across the U.S. have been filmed. The master reels must be located and new (silver negative) scanning copies must be made, but film is generally available. The physical scanning can be located practically anywhere since it is easy to ship small microfilm boxes to a remote location. The scanning process itself is quicker, easier, and cheaper with film than with paper.

The major drawback to using microfilm is that film images can be quite poor, and image quality can vary widely even on the same reel. This second point makes it difficult to scan some reels because scanning settings should be changed or adjusted in the middle of a reel. Digital image quality is directly dependent on the quality of the image on the film, on how well the microfilm operator performed, perhaps decades ago, when the film was made.

**SCANNING, ARTICLE-LEVEL, AND METADATA**

The digitization process starts with a newspaper page, either an original or on microfilm. The first step is to scan the source materials, which create 400 dpi, grayscale, digital images as TIFFs. Software crops the images (to eliminate edges on the frame that are not part of the page), deskews them (to orient the page directly up and down), and despeckles them (to remove small blurs or dots picked up by the scanner).



The homepage of Utah Digital Newspapers project

Each page goes through an article “zoning” process where human beings identify and classify each article on the page as news, an advertisement, or birth, death, or marriage announcement. (Some digitization processes use software to segment the newspaper page into article zones. This automated solution, however, has yet to prove to be reliable enough for us to incorporate into our process. We much prefer the higher quality and reliability that comes from our manual process.)

Categorizing article types allows for narrower searches, e.g., searching birth, wedding, or death announcements. Two people separately transcribe the masthead and article headlines and subheadings. This ensures that headings and subheadings are at a 99.5% or greater accuracy rate, which enhances search accuracy above and beyond normal, uncorrected OCR. (Our statement of work with our supplier requires all manual data entry be double-keyed and reconciled.)

Segmenting pages into separate articles has several advantages: Articles are presented in search results, article images are presented for viewing, and OCR is improved because there are more consistent fonts within an article and hyphenated words are more easily conjoined. We find these reasons compelling enough to justify the additional expense of segmentation. Indeed, 67% of UDN users rate the search accuracy as “good” or “excellent.”

The simple solution to organizing newspaper metadata is to follow the NDNP specifications. They are rapidly becoming the industry standard, and any viable digitization processor should be able to deliver newspaper files in this

## SUGGESTED READINGS

- For a similar article to this tailored for historians, see John Herbert and Karen Estlund, "Creating Citizen Historians," *Western Historical Quarterly*, Vol. XXXIX, No. 3, August 2008 (forthcoming).
- For materials on the early history of UDN, title selection, an overview of the digitization process, website organization, and bitonal versus grayscale images, see John Herbert and Kenning Arlitsch, "digitalnewspapers.org: The Digital Newspapers Program at the University of Utah," *The Serials Librarian*, Vol. 47, Nos. 1–2, 2004.
- For materials on UDN grant history, originals versus microfilm, file storage, distributed content, optical character recognition, multiple word options, OCR and search accuracy, and aggregated searching, see Kenning Arlitsch and John Herbert, "Microfilm, Paper, and OCR: Issues in Newspaper Digitization," *Microform and Imaging Review*, Vol. 33, No. 2, spring 2004.
- For a presentation style document (i.e., PowerPoint slides) on how to launch a digital newspaper project, see "Digital Newspapers Project Handbook" on our digital newspapers website, [www.lib.utah.edu/digital/unews/pdf/project\\_handbook\\_4.pdf](http://www.lib.utah.edu/digital/unews/pdf/project_handbook_4.pdf).

format. Even if NDNF funding is not being used, following its metadata guidelines and technical specs will better prepare you for the time when that funding is solicited.

### OPTICAL CHARACTER RECOGNITION

OCR is the centerpiece of creating full text from digital images. Today there are literally hundreds of commercially available OCR packages with wide ranges in sophistication and price. Digitizing historic newspapers is much more complex and difficult than generating text from many other document types. Some common problems are deteriorated originals, unusual fonts, faded printing, shaded backgrounds, fragmented letters, touching/overlapping letters, skewed text, curved lines (which is very common in bound volumes), and bleed-through. In fact, running OCR against a grayscale image (rather than a black-and-white image) can actually reduce OCR accuracy where bleed-through has occurred. Consequently, nothing but the most robust OCR software should be used for historic newspapers.

One method of measuring OCR effectiveness is how accurately it determines which words are on the printed page. This is normally expressed as the percentage of words on the page that are accurately "read" by the software. The software reads words letter-by-letter (character-by-character), so OCR accuracy is often measured as the percentage of characters accurately "read." It is important to distinguish between character accuracy and word accuracy: Word accuracy is by definition significantly lower than letter accuracy because it combines the probabilities of all the letters in the word. For example, OCR accuracy at the letter-level for a document may be 98%. But computing the accuracy of a five-letter word in that same document is done by taking 0.98 to the fifth power (the joint probability of five letters), which is 90.4%. This distinction between letter and word accuracy is critical when analyzing OCR accuracy.

One other consideration with newspaper OCR is that articles often have a great deal of redundancy. The same impor-

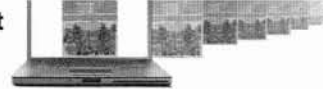
tant search keyword, such as a last name or the name of a city, often appears more than once in the same article. In these cases, word accuracy need not be 100% for a successful search.

### GENERAL DESCRIPTION OF OCR SOFTWARE

Included in the processing services offered by our service provider, iArchives, Inc. of Lindon, Utah, is OCR. The company's OCR software is state-of-the-art, proprietary, and patented. Consequently, what follows is a fairly generic description of how OCR software operates.

iArchives' OCR framework uses multiple engines, each of which runs with a different orientation. For instance, one engine may work better with lighter images and thinner fonts, while another may work better with darker images and bolder fonts. These orientations are needed because images span a wide range of quality from article to article and page to page.

Each engine inspects each image pixel by pixel, looking carefully at contiguous dark pixels, determining their overall shapes, and comparing those shapes to known letters in many different fonts. It also closely examines the adjacent white areas. One of the most important decisions the software makes is, given the size, shape, and location of a white space, is it 1) part of a letter that is fragmented, 2) the space between letters, or 3) the space between words? Various algorithms are run to help answer these questions, including growing and shrinking potential letter fragments to see if dark areas can be connected to form letters; whitening the background and darkening the text to improve the contrast; and changing near-white to white and near-black to black. In the end, the software assembles the contiguous/connected dark pixels into a letter, sometimes noting that more than one letter is a possibility. Letters within a small distance of each other are assembled into a "node," which is what the engine believes is a word.



Finally, because the engine may have multiple possibilities for a particular letter, it may as a result have multiple possibilities for the associated word that the letter is in. For example, if the engine can't decide about the third letter in the node "be\*t," possible word options include beat, beet, bent, and best. The engine will rank order the possible words from best to worst option, and then, depending on the setup parameters, provide the requisite number of words for the text. These words are accumulated into the text file for the article (included with the text and the x- and y-coordinates of each node's location within the image, which are used to highlight the word when a successful search for it is done).

#### THE DILEMMA OF MULTIPLE WORD OPTIONS

When searches are performed on the entire UDN collection, the search-hit limit of 20,000 can easily be reached, especially if a single and somewhat common word is used in the search. One method of limiting a search, of course, is to search on more than one word. As is quite popular, a user will search on a first and last name together, instead of the last name only. For example, searching the collection for *cassidy* results in 1,932 hits, but when that is limited by searching for *butch cassidy*, 201 hits result. Additionally, if an exact-phrase search is used to further limit the search, only 161 hits result. This is where things get interesting.

On page 2 of the *Eastern Utah Advocate*, Nov. 14, 1901, there is an article titled "Eastern and Southern Utah." On the printed page is the phrase "George Parker alias Butch Cassidy"—the OCR generated the following corresponding text: "george parker alias alfas butch dutch cassidy." The OCR had some difficulty with the words "alias" and "Butch," putting two options in for each. In particular, the two options for "Butch" (butch and dutch) present an intriguing problem. The three words "butch dutch cassidy" as they appear in the text preclude a successful search on the exact phrase "butch cassidy" because "dutch," as the second word option for "butch," appears in between.

So the dilemma is this: When the OCR generates more than one word for a particular node, we enhance our ability to successfully search on that single word. In the example above, both a search for *butch* and for *dutch* will generate a hit on the article. At the same time, almost paradoxically, a multiple-word option reduces the chances of a successful exact phrase search. In fact, an exact-phrase search for "butch dutch cassidy" reveals six hits, so the identical OCR result occurs five other times in the database.

After the initial text is generated by the OCR, it is filtered through a number of different dictionaries to ensure that only valid words are in the final text. We feel surnames are particularly important given the high genealogical use of the collection. iArchives utilizes a 2-million item dictionary containing all English words, common foreign language words, surnames, and place names. Additionally, we augmented the place names dictionary with a set of Utah place names provided by BYU.

#### IMAGE QUALITY

Up to now, we have stored an image of each article in a separate file, as well as an image of the full page. These images are black-and-white PDF files with embedded text. We occasionally use grayscale PDF files for papers if they require a higher level of quality and detail to read the text, such as the early *Salt Lake Tribune*, which has up to 10 columns on a page.

As technology develops and more users switch to higher speed internet connections, we will advance to higher quality and more user-friendly images. This presents a significant trade-off between image quality and speed of download. We are transitioning to grayscale images in the JPEG 2000 format, which allows the user further zooming capabilities, to "clip" a section of a page, and to highlight an article on the page image.

#### FUNDRAISING CONSIDERATIONS

As the Utah Digital Newspaper program continues to develop into the future, it will incorporate technological innovations to aid in searching or browsing. The most frequent suggestion from our users is to "add more content," and we intend to oblige. We hope to broaden our user community to more researchers and the K-12 audience. In order to provide continued long-term growth, our goal is to become financially self-sustaining.

Being exclusively a soft-money program means we constantly need to raise money to expand the collection. To date, that has amounted to an impressive \$2.5 million, but our needs continue to exceed our resources.

We have had tremendous success in our first 5 years, digitizing 600,000 pages. But with our long-term goal exceeding 8 million pages, we have only scratched the surface. We still need many millions of additional dollars to reach our fundraising goal and to begin to consider our collection complete.

#### Endnotes

1. Frank R. Jenkins, Thomas A. Nartker, and Stephen V. Rice, "Testing OCR Accuracy," *Inform*, Vol. 10, No. 8, September 1996, pp. 21-25.
2. Thomas A. Nartker, Stephen V. Rice, and Frank R. Jenkins, "OCR Accuracy," *Inform*, Vol. 9, July 1995, pp. 38-40.
3. Kazem Taghva, Julie Borsack, Allen Condit, and Srinivas Erva, "The Effects of Noisy Data on Text Retrieval," *Journal of the American Society for Information Science*, Vol. 45, No. 1, January 1994, pp. 50-58.

*John Herbert* ([john.herbert@utah.edu](mailto:john.herbert@utah.edu)) is head of digital technologies, J. Willard Marriott Library, University of Utah, and *Karen Estlund* ([kestlund@uoregon.edu](mailto:kestlund@uoregon.edu)) is digital collections coordinator, Knight Library, University of Oregon (formerly with The University of Utah).  
Comments? Send email to the editor ([marydee@xmission.com](mailto:marydee@xmission.com)).

