

Using Chronicling America's Images to Explore Digitized Newspapers and Imagine Alternative Futures

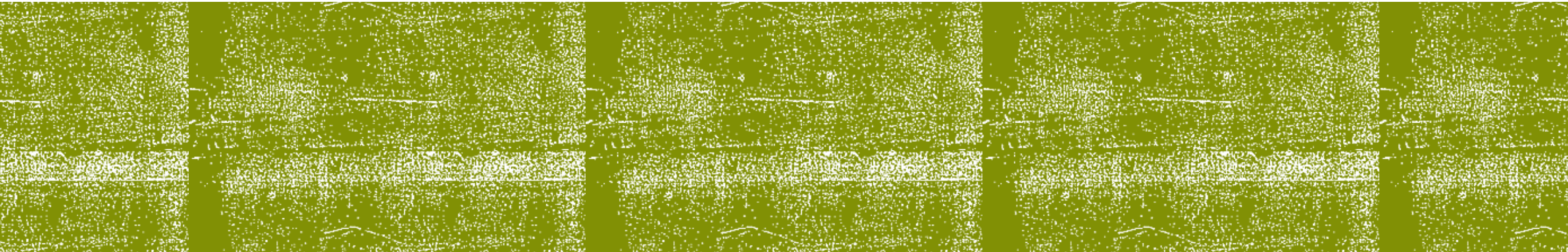
Elizabeth Lorang, Leen-Kiat Soh, Yi Liu,
Chulwoo Pack, and Delaram Rahimi



University of Nebraska–Lincoln
projectaida.org



Researchers and developers might do far more with digital images at **all stages of digitization and use**, and that attention to the digital images will yield greater understanding across a range of domains

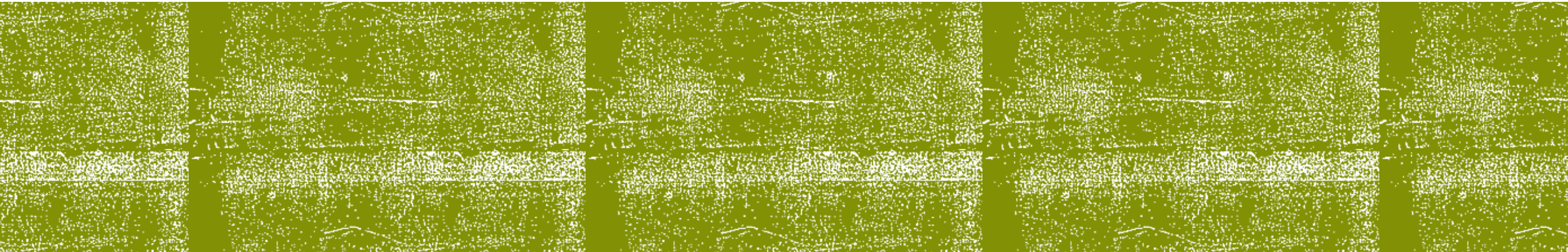




Much to learn about the **materials** themselves as well as about **process**: (1) values we bring to digitization and those that get enacted through digitization, and (2) how current work extends from those values in maintaining the original items

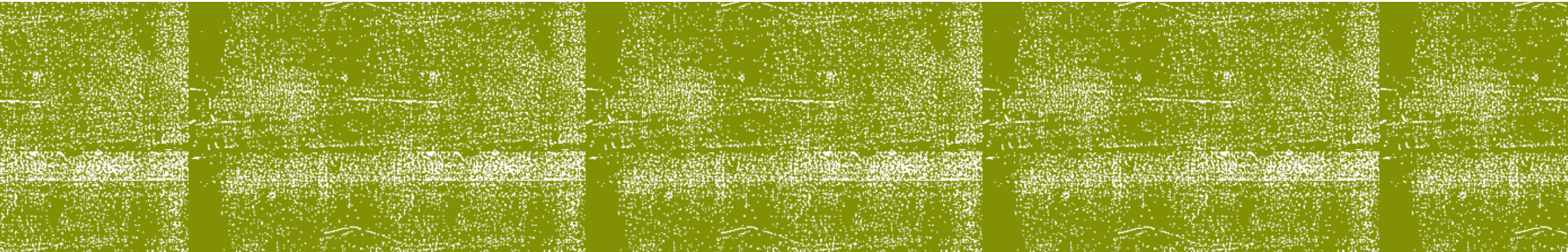


Studying these processes can function as a form of social and cultural history: how society and culture are captured in collecting and in technology



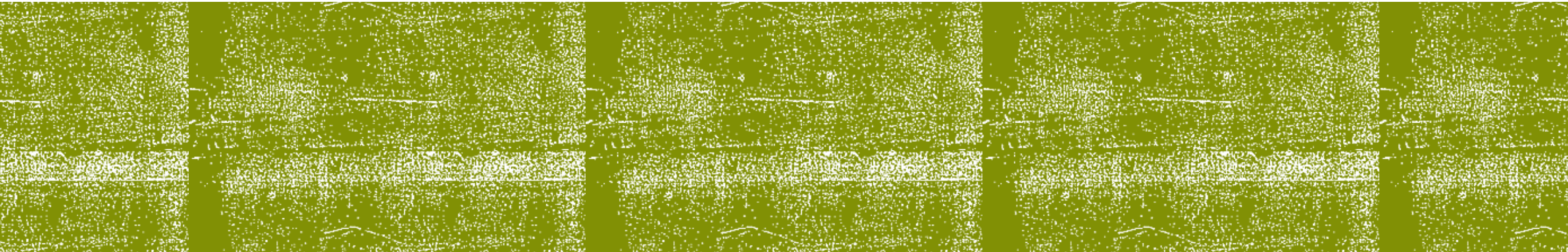


Exploration and study of digital images of historic materials as a mode of asking questions about the **materials of many forms**: physical originals, microform duplications, digital copies, represented as images, text, and metadata





What might we learn about digital collections of cultural heritage materials, and how might we augment use and access of these collections, if we focus on the digital images being created as librarians, archivists, museum professionals, and others are digitizing cultural heritage materials?





POETRY

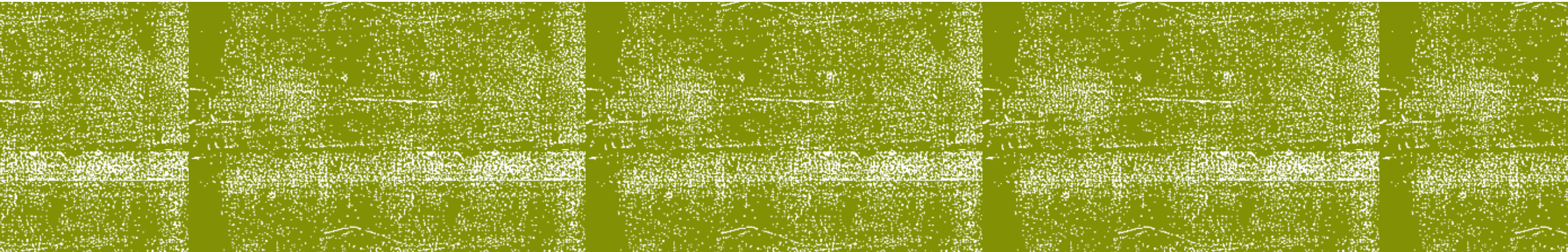
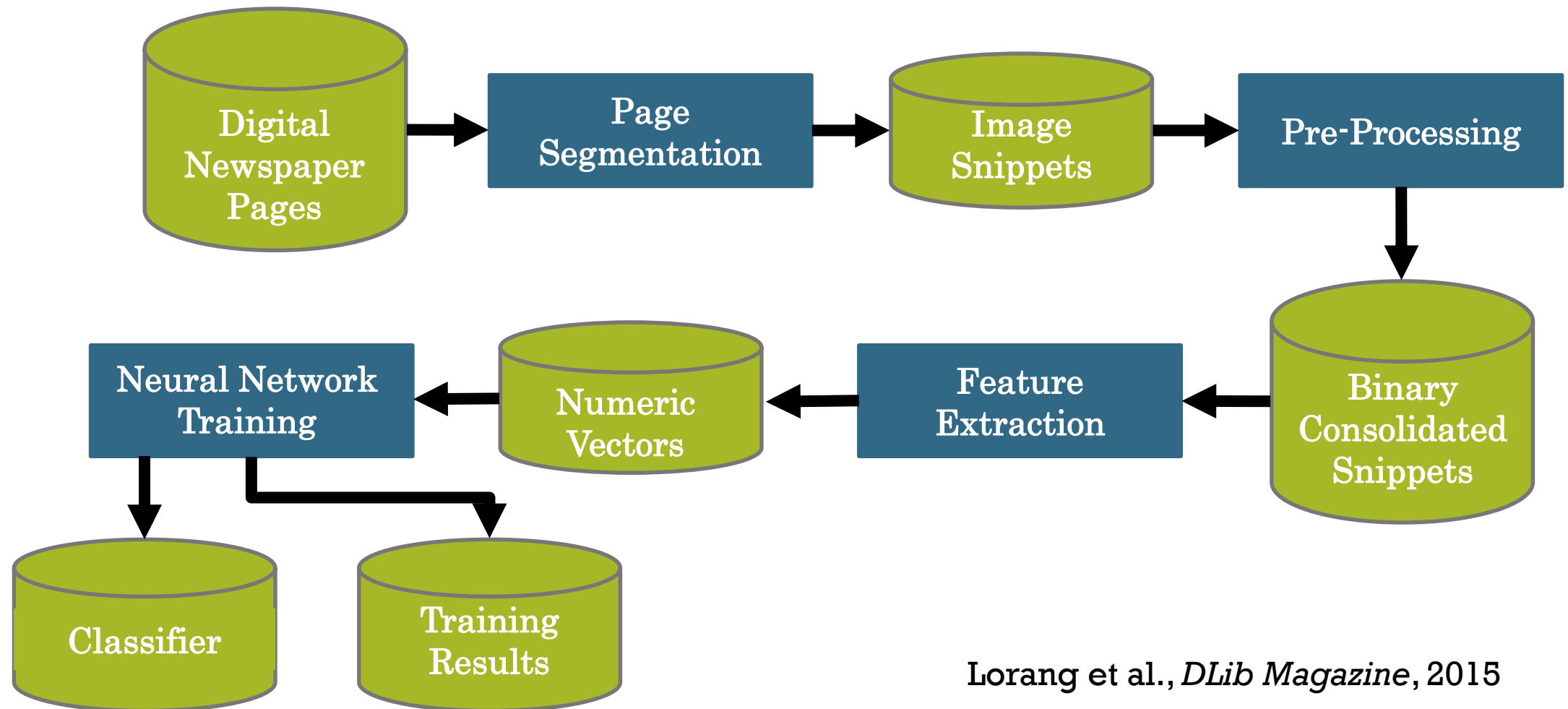


IMAGE PROCESSING

Generate data about visual features from the newspaper pages and then use those extracted features within a computational system called an artificial neural network



Lorang et al., *DLib Magazine*, 2015

CASE STUDY

Test the effectiveness of our first-generation approach on page images from Chronicling America from 1836-1840 (8,000 pages, downloaded in 2016)

The results lagged well behind what we saw in our smaller proof-of-concept project

NOISE

WASHINGTON HALL,
University of Virginia.

At a called meeting of the Washington Society this evening, Mr. James L. Orr, arose and said, Mr. President, it becomes my painful duty to announce to the Society the death of one of our honorary members, *Thomas Butler Bird*, of South Carolina. When the melancholly intelligence first reached this place, some faint shadow of hope, as to its truth, prevented our giving entire credence to the tragical affair.— But it is now too sadly confirmed, and our much esteemed friend sleeps in death's icy embrace. He has been cut off in the spring time of his existence, and we are left to weep over the many generous qualities of his nature—the bud was just opening—its promised fragrance was adding new charms to its loveliness—but alas! it has been thus early nipped by an untimely frost, and consigned to wither and decay.

"All that's bright must fade,
The brightest still the fleetest;
All that's sweet was made.
But to be lost, when sweetest."

When we reflect that he was distinguished alike for his benevolent spirit, a nobleness of heart, and a superiority of talents, the sympathetic tear starts to swim the eye and moisten the cheek, on account of his unhappy fate. I shall attempt to pronounce no eulogy on his character, but the sorrowed countenances

good quality

Save when some quick emotion
The warm blood strongly sent
To revel in her olive cheek,
So richly eloquent.

I said consumption smote her,
And the healer's art was vain;
But she was an Indian maiden,
So none deplored her pain;—
None, save that widow'd mother,
Who now, by her open tomb,
Is writhing like the smitten wretch
Whom judgment marks for doom.

Alas! that lowly cabin,
That couch beside the wall,
That seat beneath the mignon vine,
They're lone and empty all.
What hand shall pluck the fall, green corn,
That ripeneth on the plain,
Since she, for whom the board was spread,
Must ne'er return again?

Rest, rest, thou Indian maiden!—
Nor let thy mummuring shade {scorn
Grieve that those pale-brow'd ones with
Thy burial-rite surveys;—
There's many a king, whose funeral
A black-rob'd realm shall see,
For whom no fear of grief is shed,
Like that which falls for thee.

Yes, rest thee, forest-maiden!
Beneath thy native tree:
The proud may boast their little day,
Then sink to dust like thee!
But there's many a one whose funeral

bleed-through

We may forget, but wilt
We mourn a lasting vacancy
No other one can fill.
Hearts that have lost no dear a friend
Must own one vacant cell,
There may be those that they can love,
But none they'll love as well.

There must be dark parental hopes
Are dashed on death's dark wave,
And she who should life's joys have shared,
Would fain have shared his grave!
Fond brothers mourn the first deep grief
They e'er were doomed to know;
Sisters, whose love may ne'er be told
Are bending beneath the blow.

And friends—who cling to each fond hope
To lose with life was gone;
And reason—O, why that word?
He had no room and
How lonely, ah! how desolate
Is childhood, and manhood, and health;
We're missing from the social ring
A part of priceless worth.

When from the gathered worshippers
Arenas the fervent prayer,
Is sung the song the bending throng,
One form no more is there,
And when the thrilling voices join,
To raise the evening song,
The harmony is not complete,
One tuneful voice is gone.

The friend, so true to many a trust,
The favored one has flown,
To find affliction's stings,
No more an answering tone.

low contrast

ANTINE, ST. JOSEPH CO.

From the Baltimore Visiter.

WE MAY BE HAPPY YET.

Ah! dearest dry those tears away,
That stain thy fading cheek;
Unbend thy lips from sorrows away,
And words of comfort speak.
Banish the past, and with me vow
Our sorrows to forget;
And be Hope's star our pilot now,—
We may be happy yet.

The care, believe me, that enshrouds
Thy cheek's once cheerful ray,
Gives me more pain than all the clouds
That darken o'er our way.
Then let thy sweet lips smile again—
Smile as when first we met,
Grief cannot always shadow them—
We may be happy yet.

Gaze on yon star so bright and clear,
Free from its cloudy chain;
Thus will our sorrow disappear,
When thou dost smile again;

occluding "blobs"

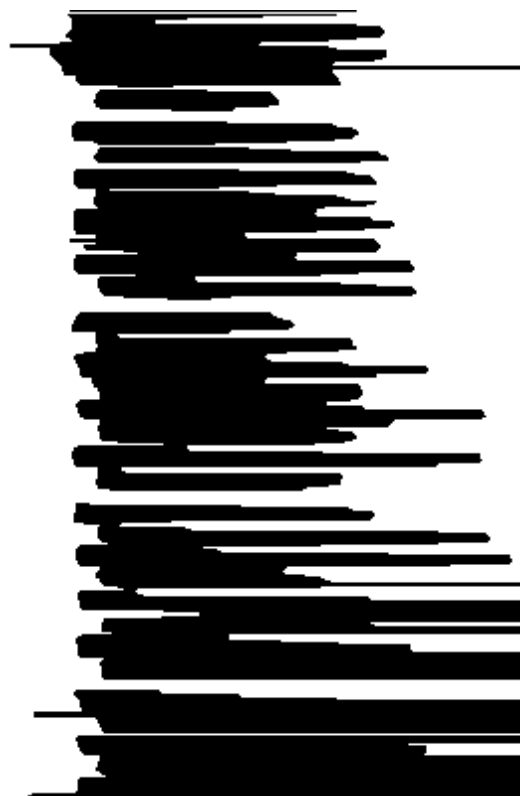
IMPROVED APPROACH

Improve our **binarization** of the image into "object" and "background" information (represented as black and white pixels, respectively) and extract visual structures with a **noise-tolerant text line segmentation** approach

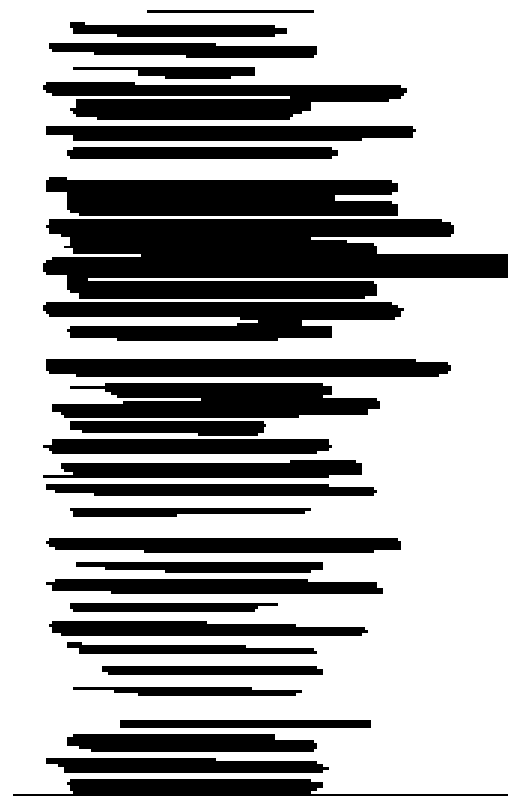
Results from Initial Approach



good quality



bleed-through



low contrast

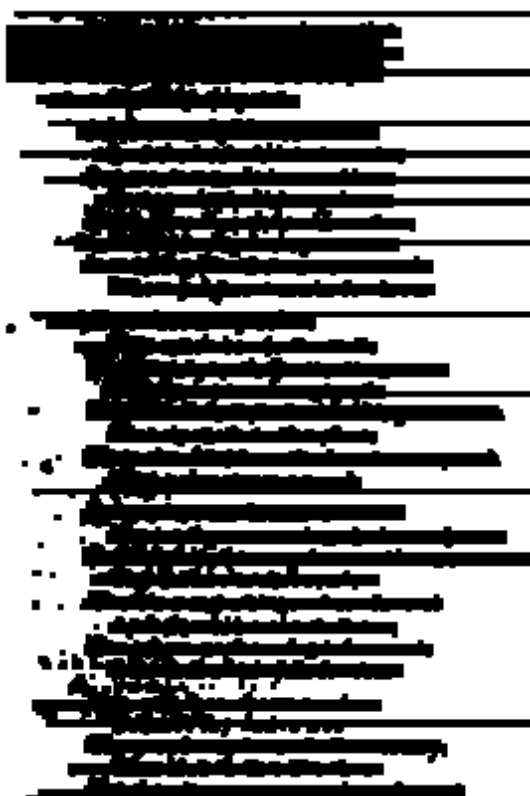


occluding "blobs"

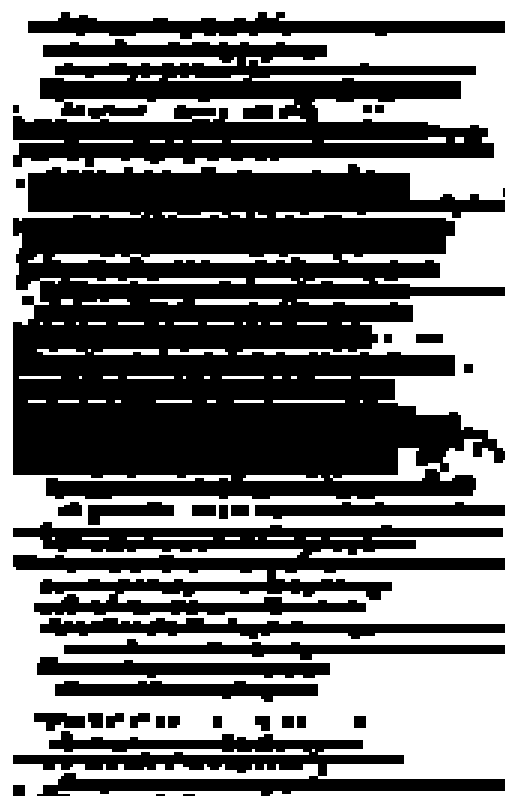
Results from Improved Approach



good quality



bleed-through

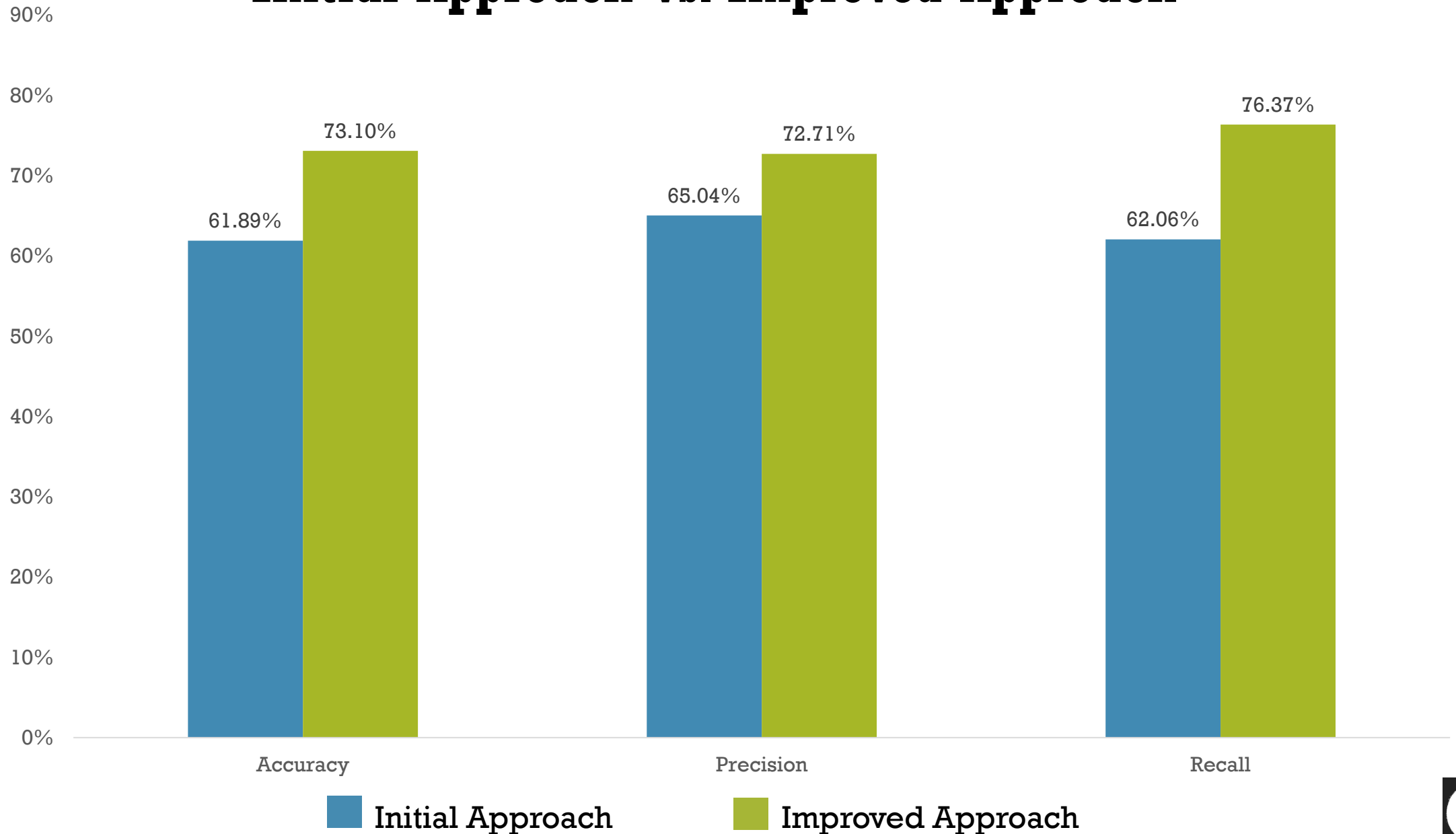


low contrast



occluding "blobs"

Initial Approach vs. Improved Approach

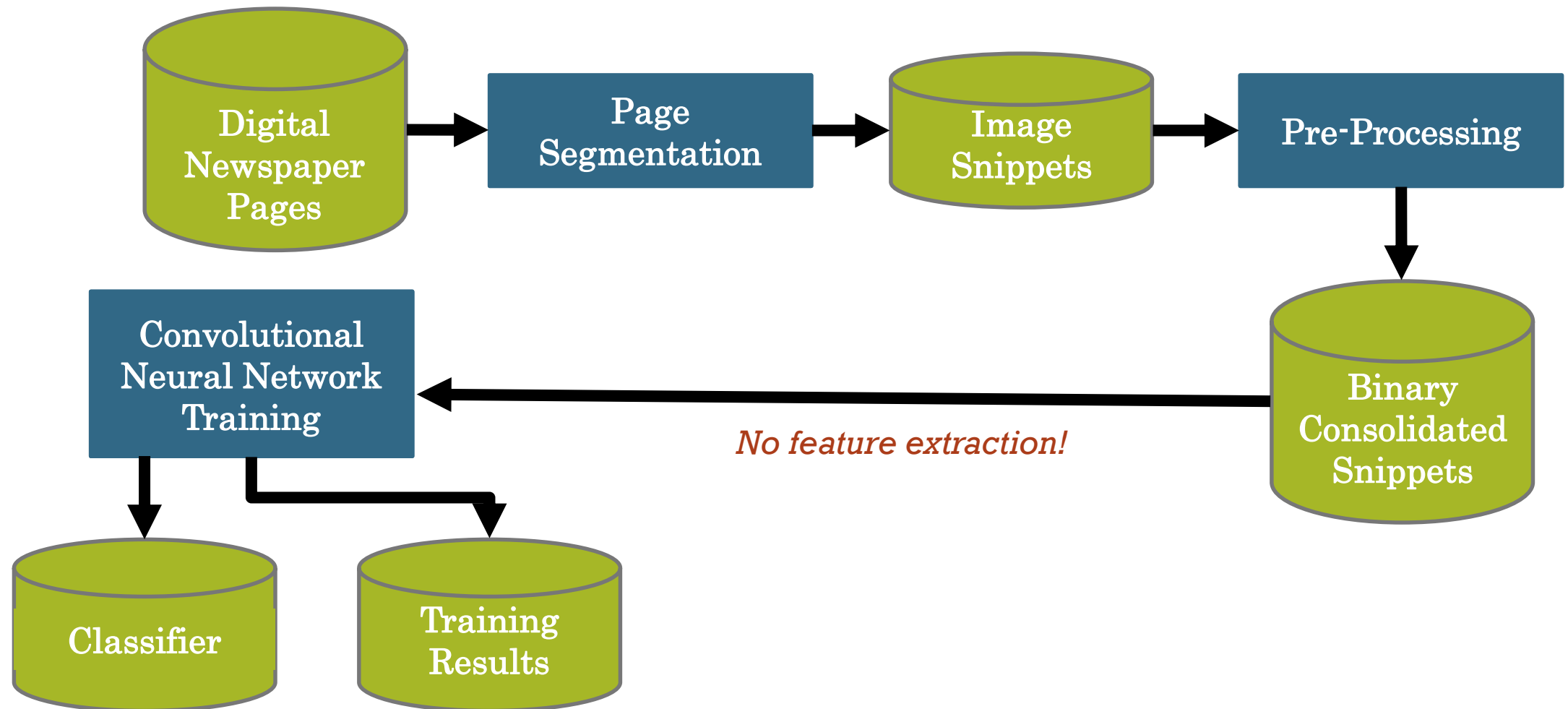


ALTERNATIVE APPROACH

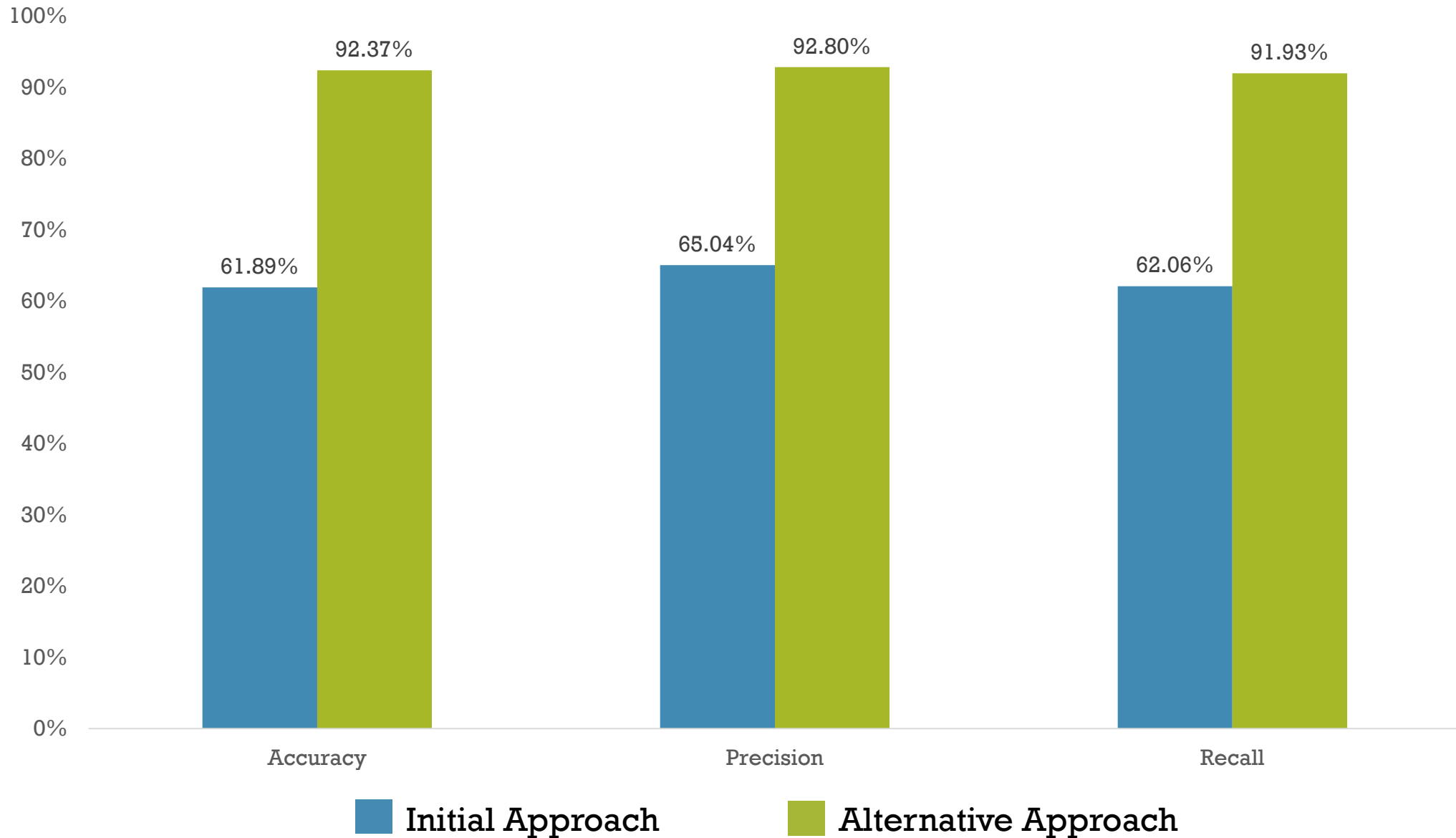
Use deep learning via a **convolutional neural network**, *without* having to perform feature extraction at all

Initial results of accuracy achieved 92.4%

- precision = 92.8%; recall = 91.9%

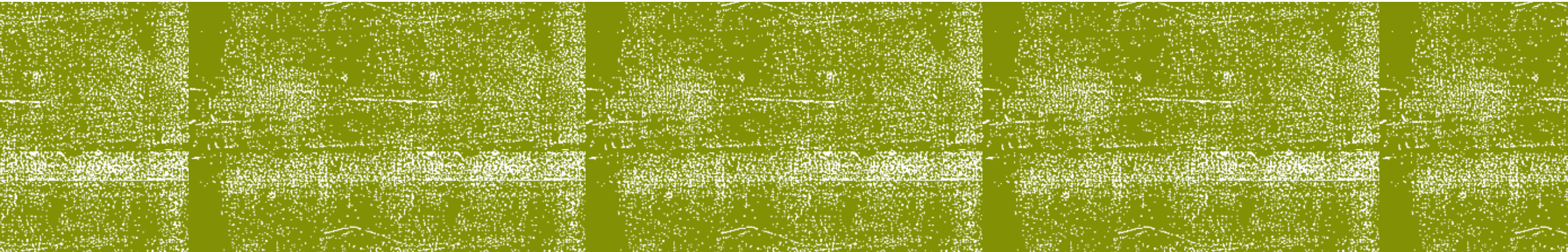


Initial Approach vs. Alternative Approach





LAYOUT ANALYSIS



LAYOUT ANALYSIS

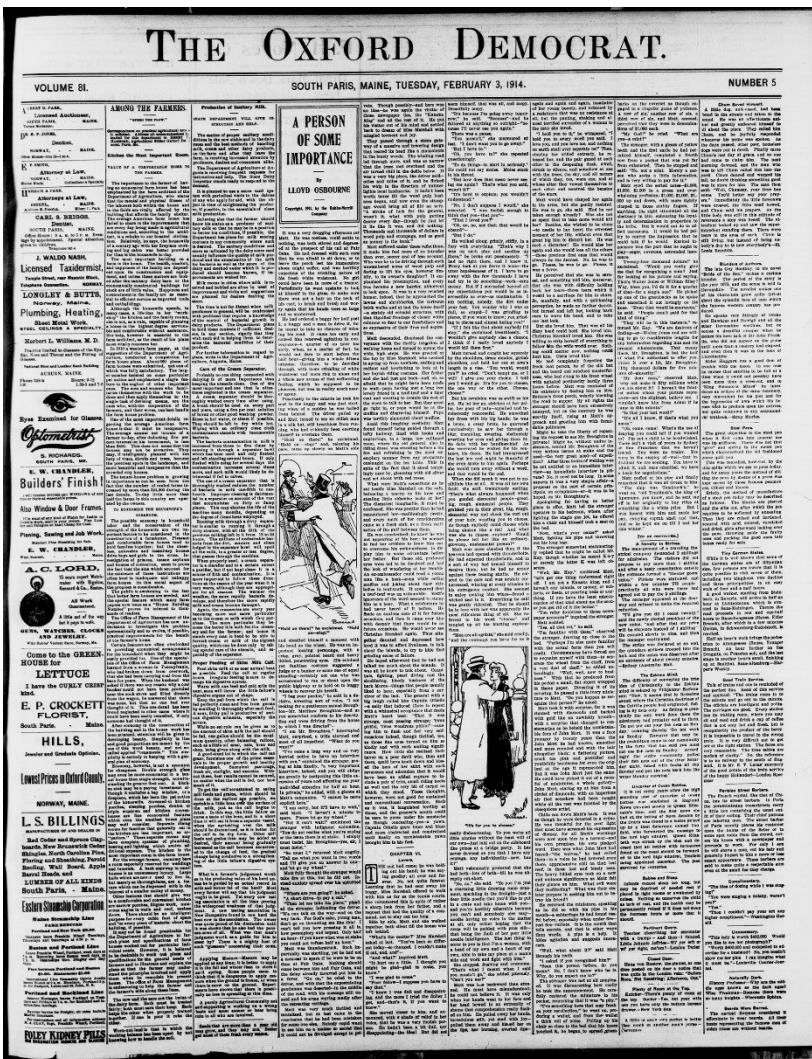
Layout analysis and zoning performed by OCR systems as a pre-processing step for text recognition

- E.g., built in to OCR web services such as ABBYY Cloud and in the open source Tesseract OCR engine

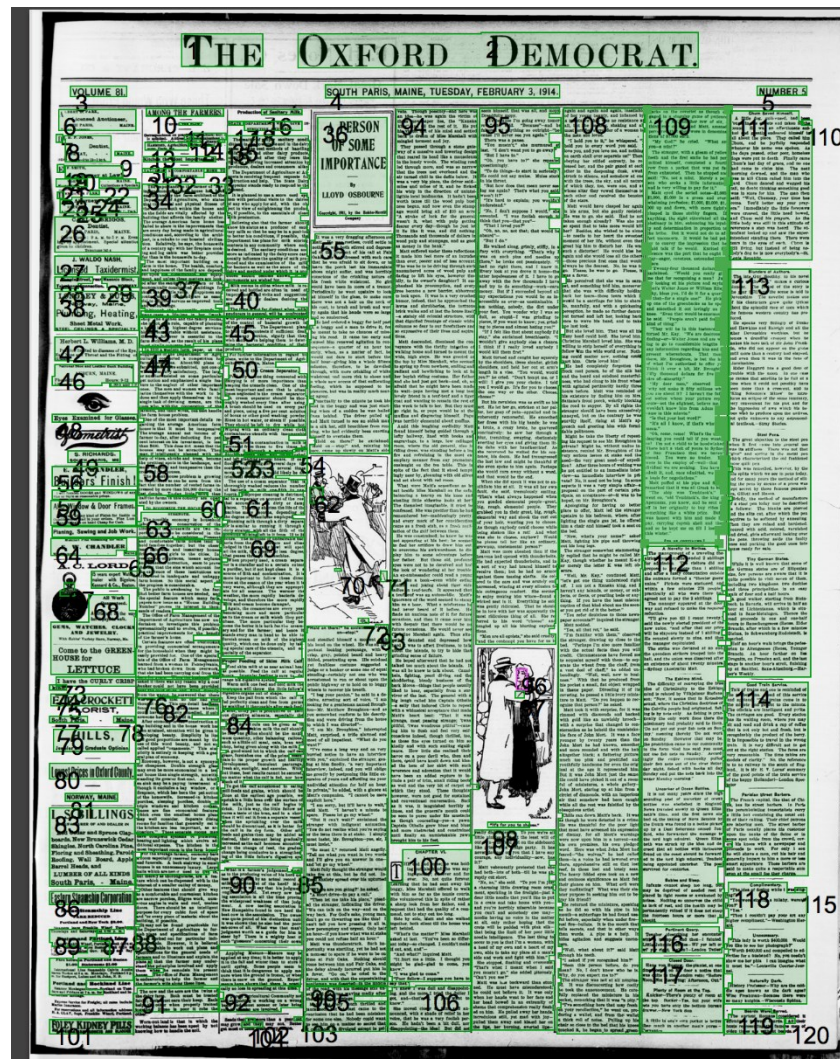
Some of this zoning information is maintained in the OCR XML files made available through Chronicling America:

- Zoning coordinates for strings, text lines, some larger text blocks, and page-level information

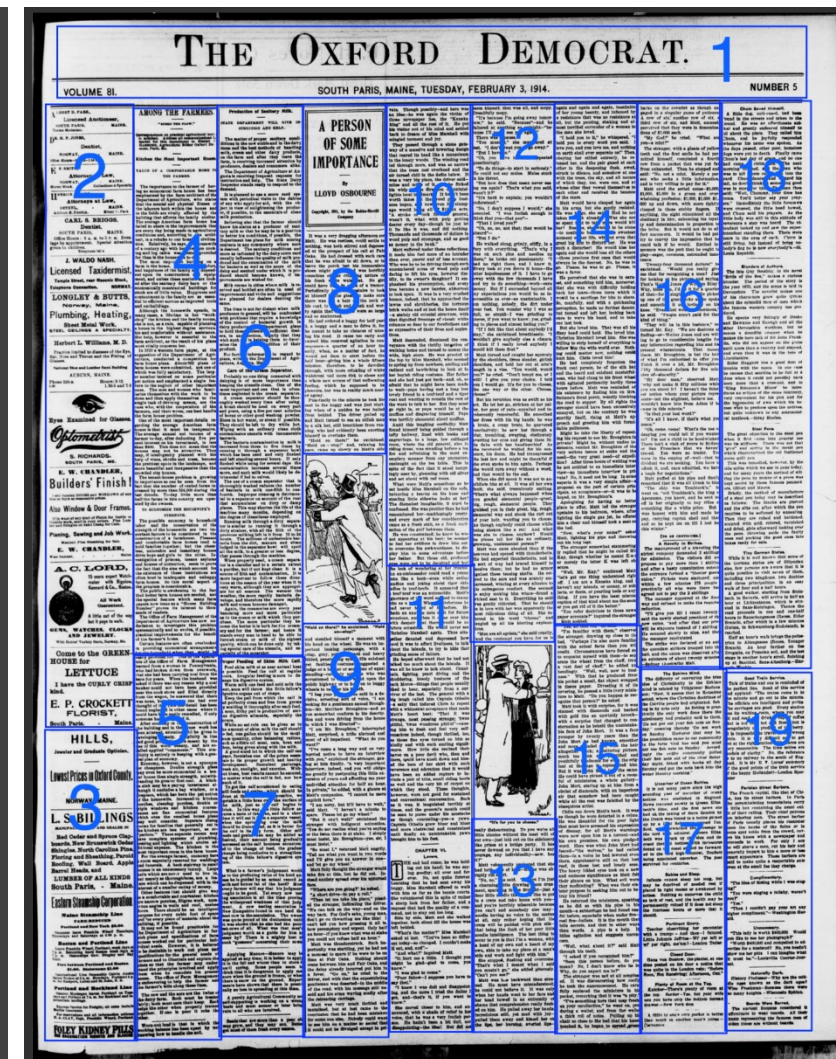
Layout analysis and zoning remaining challenging for historic newspapers and periodicals, particularly those with dense and mixed layouts



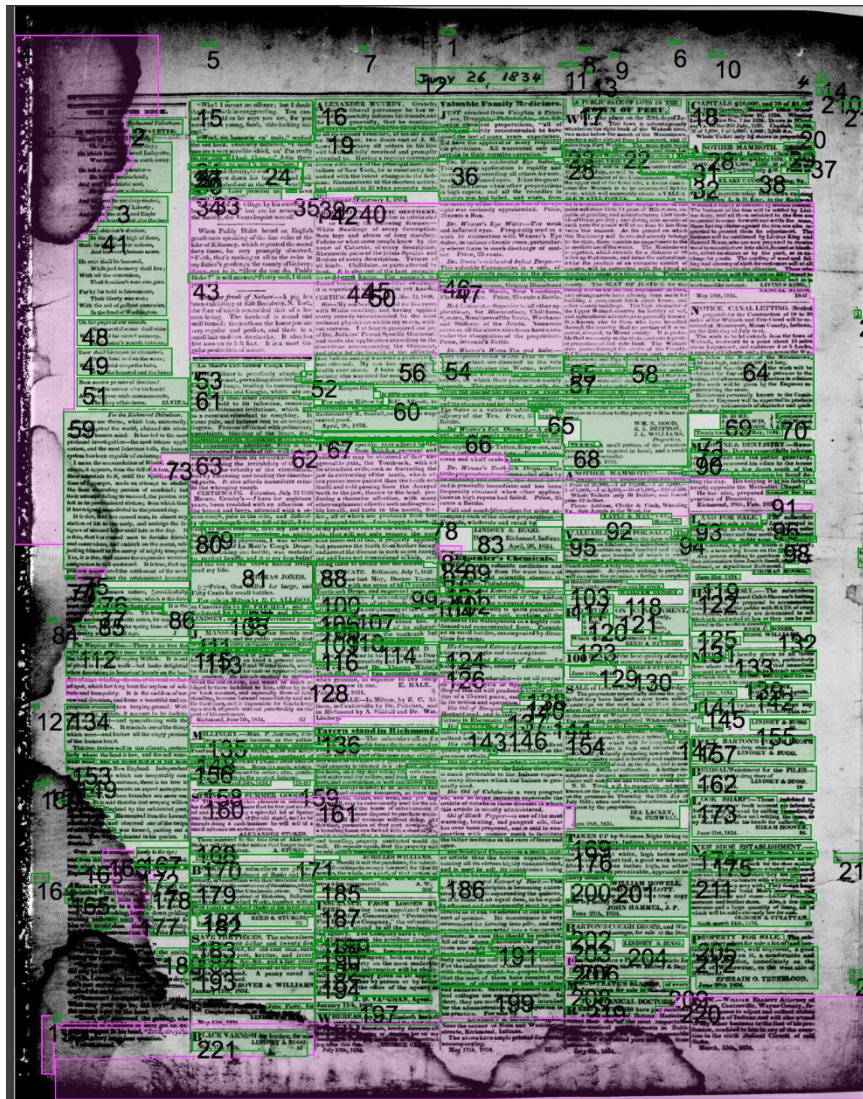
newspaper page from
Chronicling America



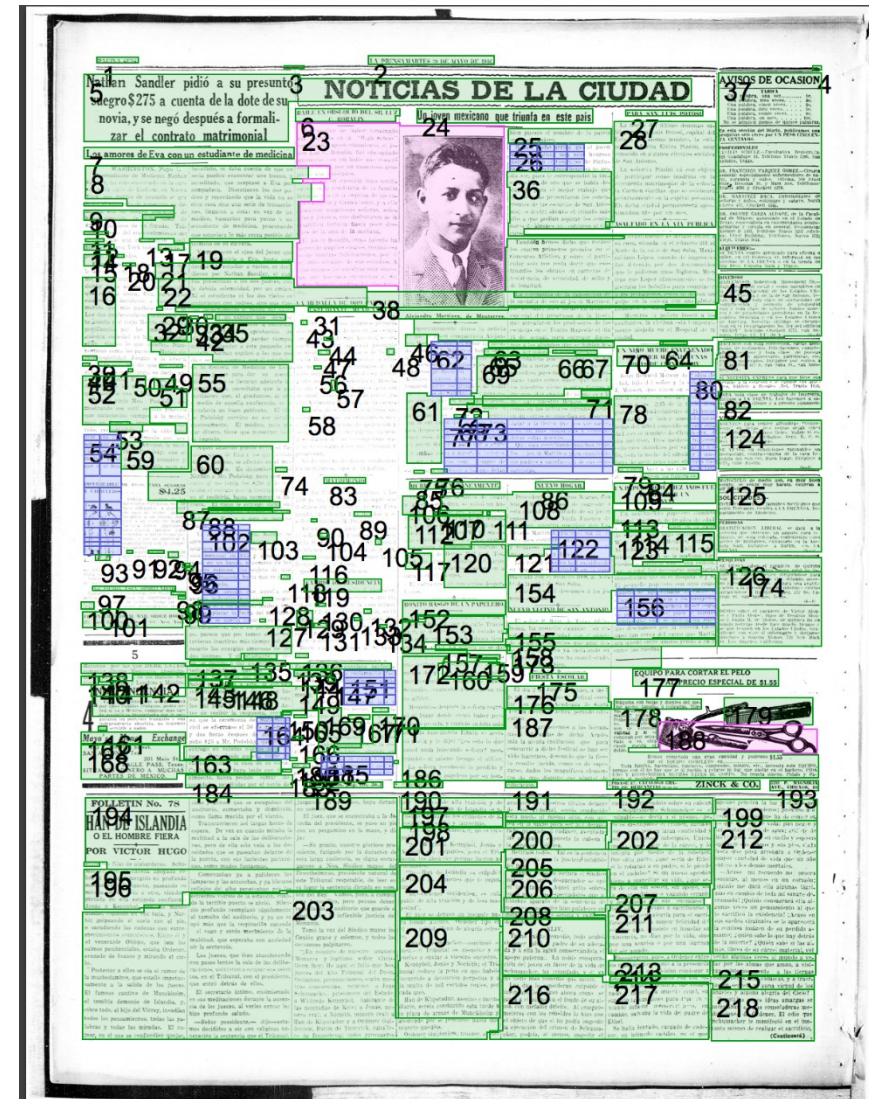
zones identified by Abbyy
FineReader



text blocks recorded in OCR XML
on Chronicling America



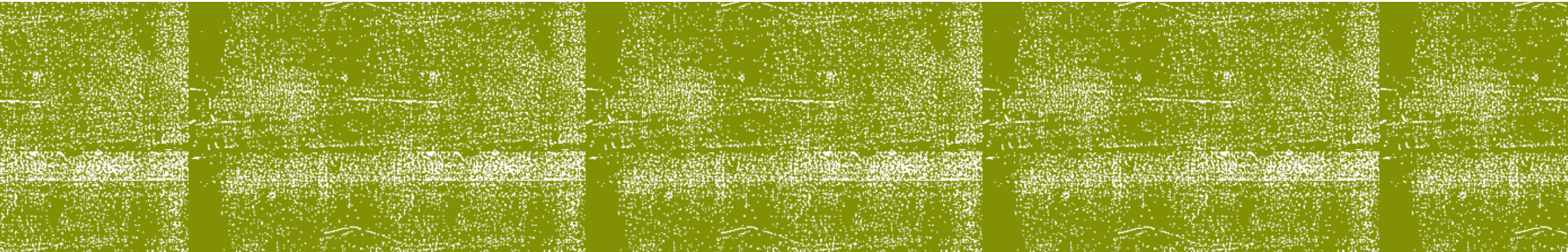
challenges to zoning:
damage & noise



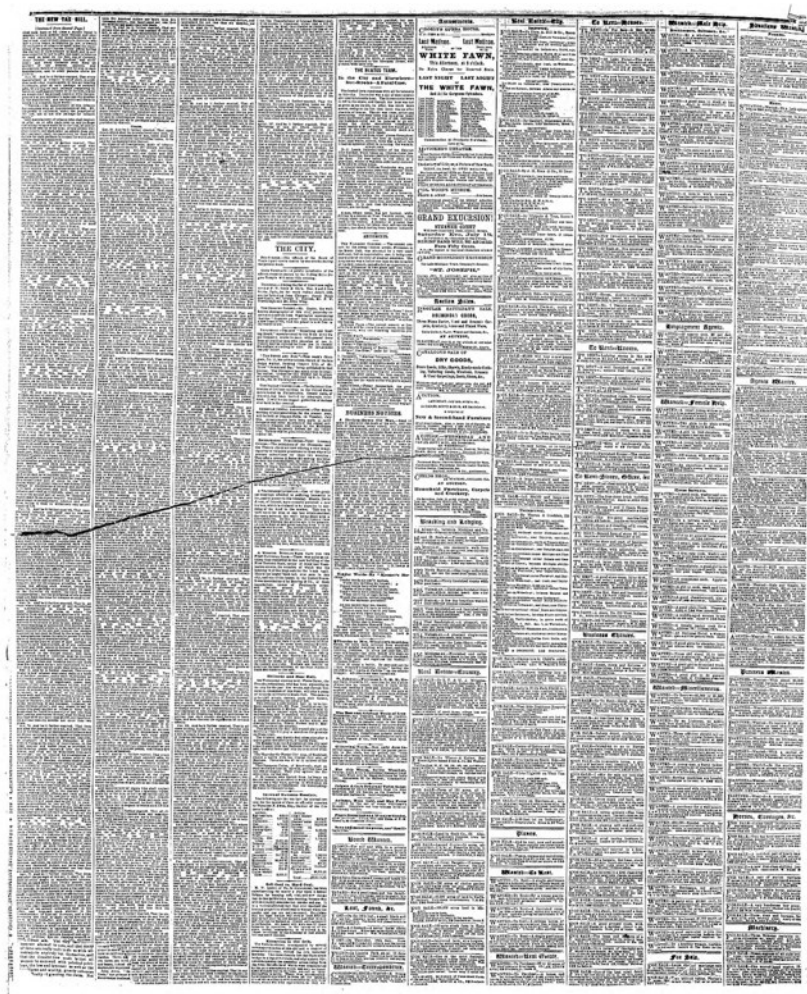
challenges to zoning: contrast
and range effect



CORPUS-WIDE ANALYSIS



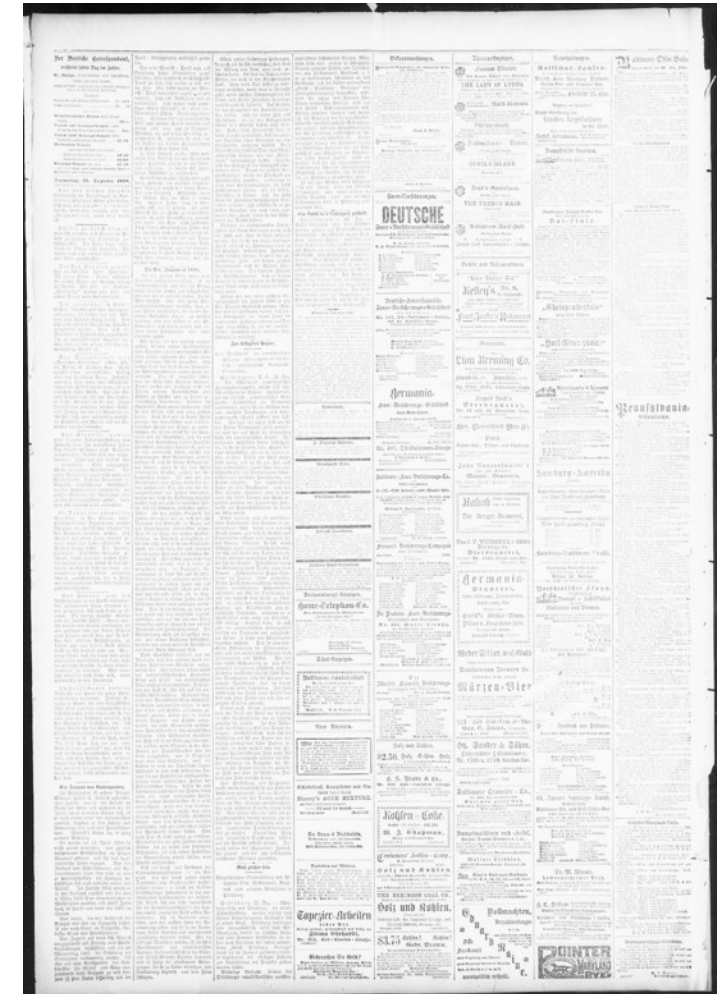
[illegible]



Newspaper page showing
high contrast



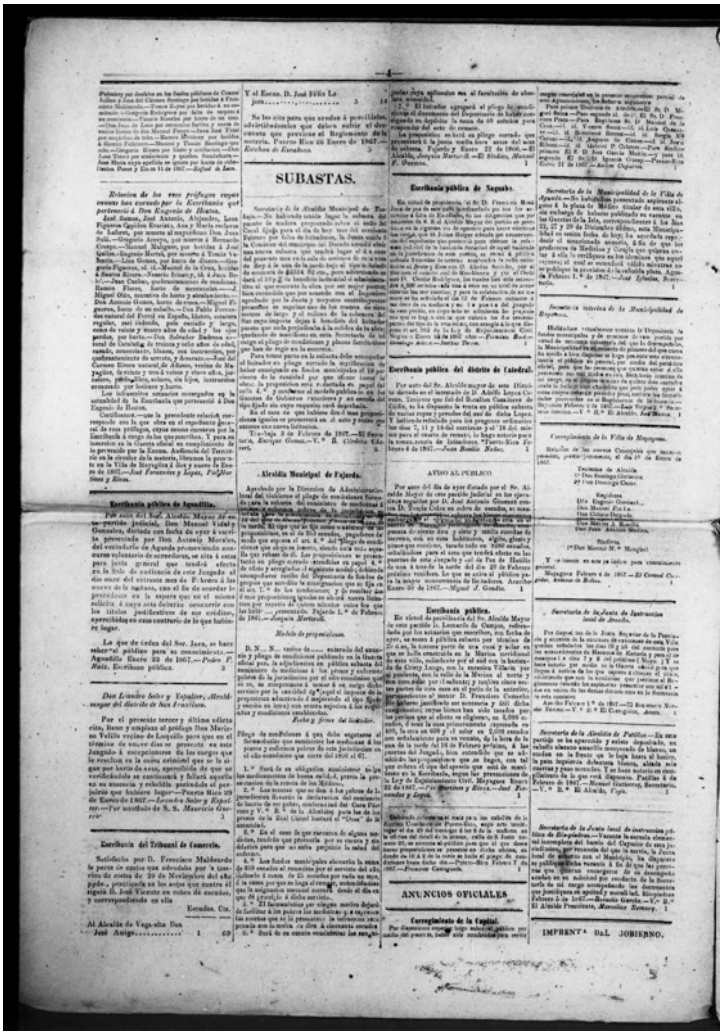
Newspaper page showing
average contrast



Newspaper page showing
low contrast



Newspaper page showing high range-effect



Newspaper page showing average range-effect



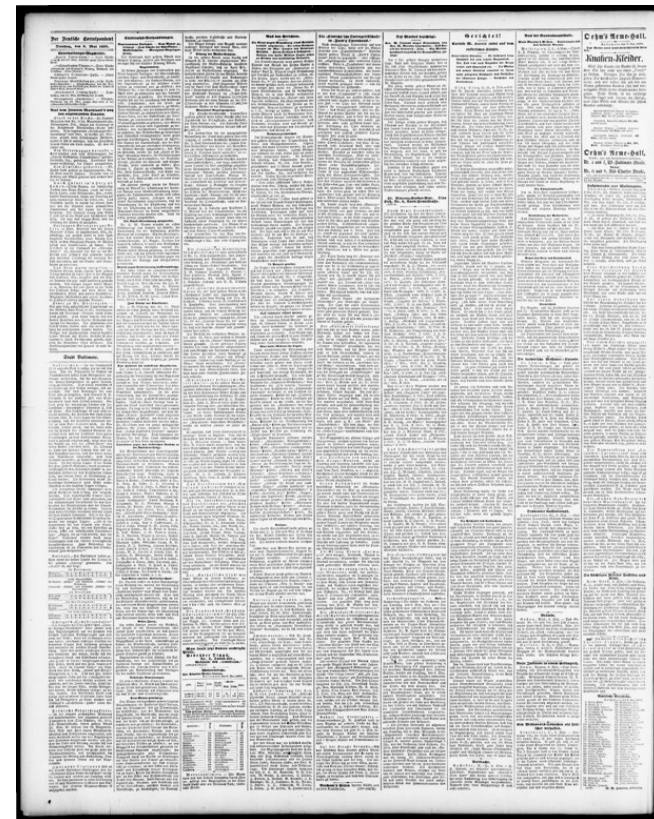
Newspaper page showing no range-effect



Newspaper page showing
severe orientation skew



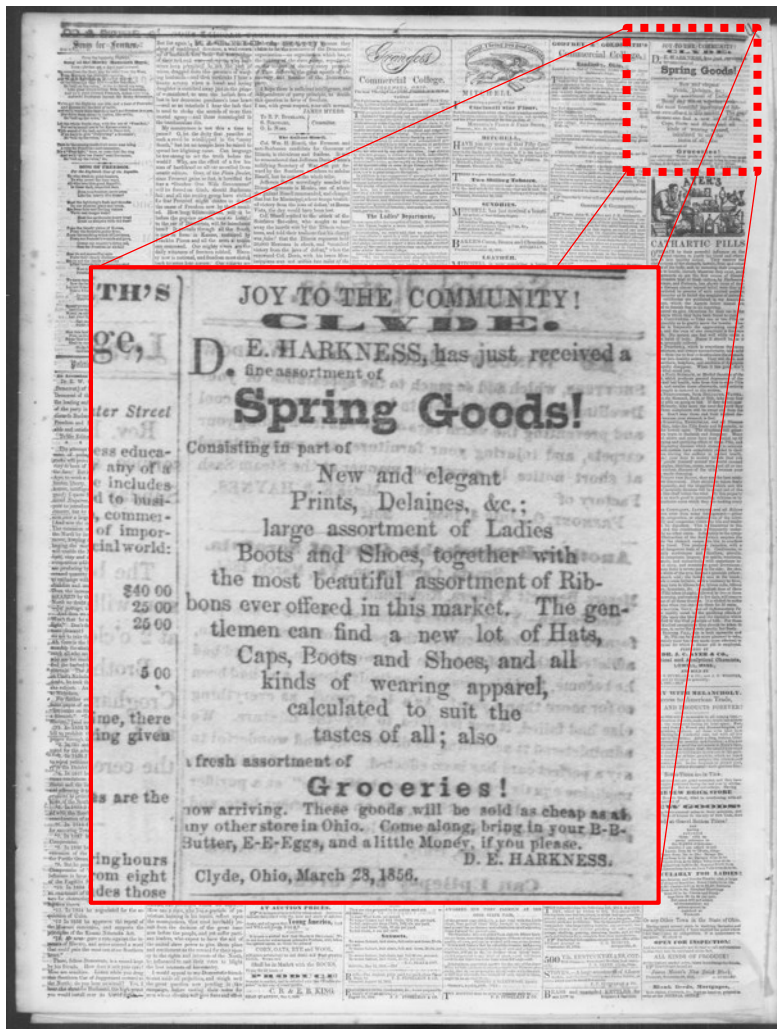
Newspaper page showing
some orientation skew



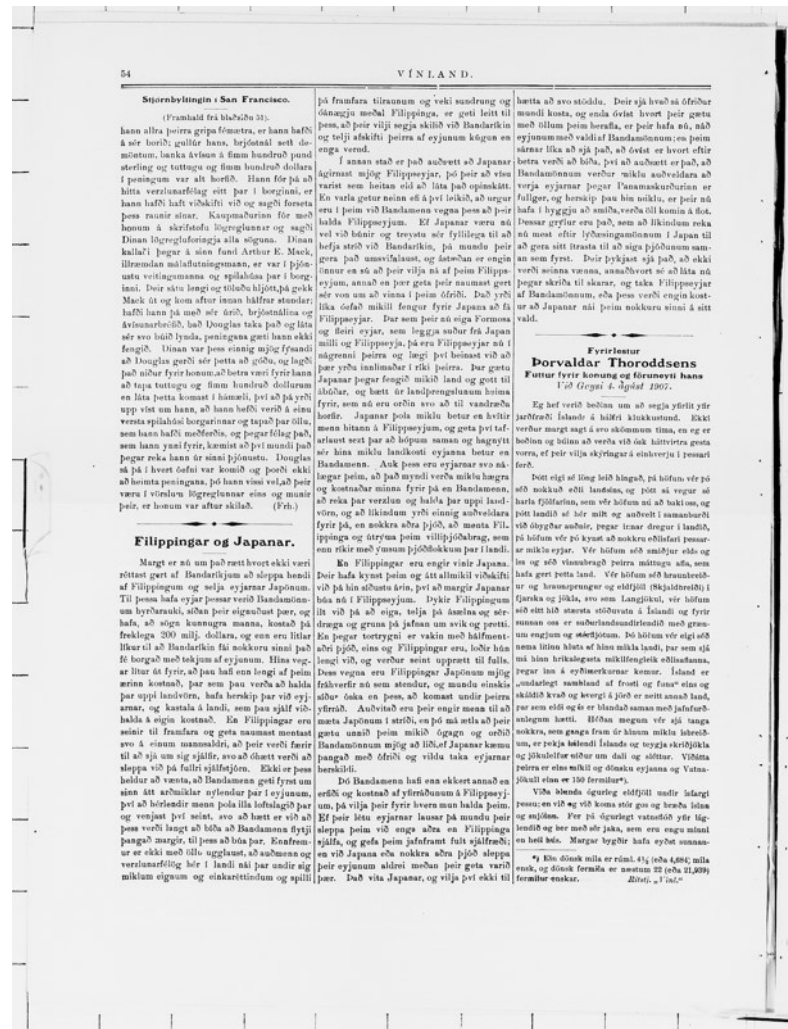
Newspaper page showing
no orientation skew



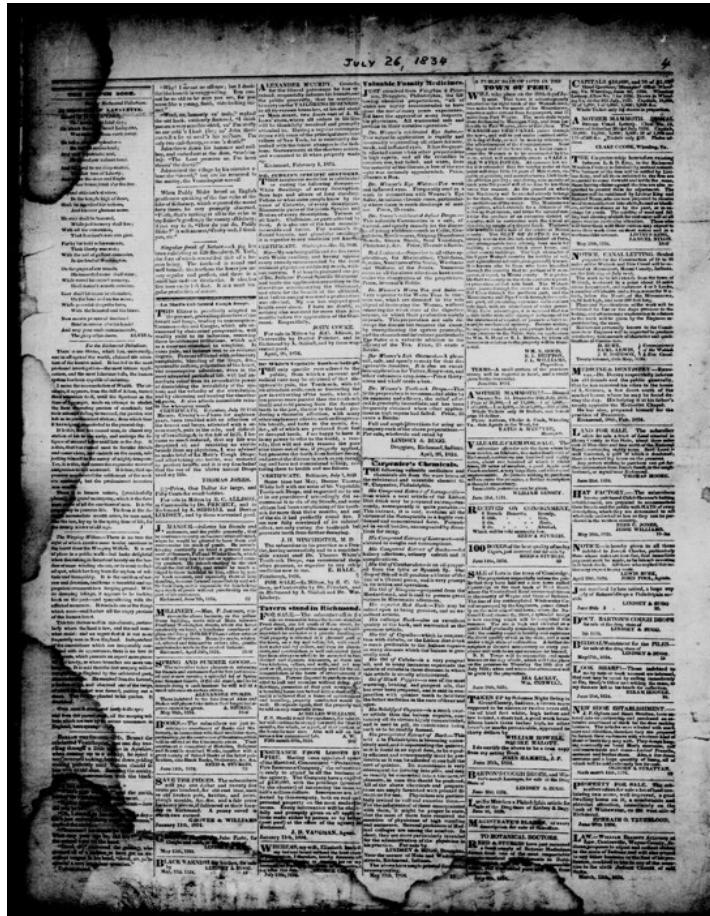
Newspaper page showing
severe bleed-through



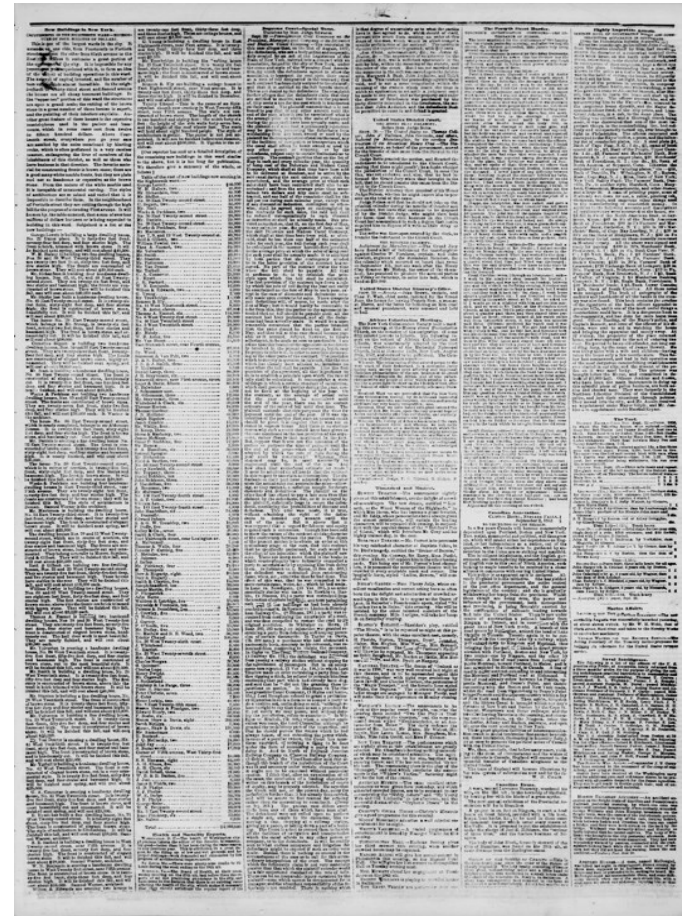
Newspaper page showing
some bleed-through



Newspaper page showing
no bleed-through



Newspaper page showing
severe noisiness



Newspaper page showing
some noisiness



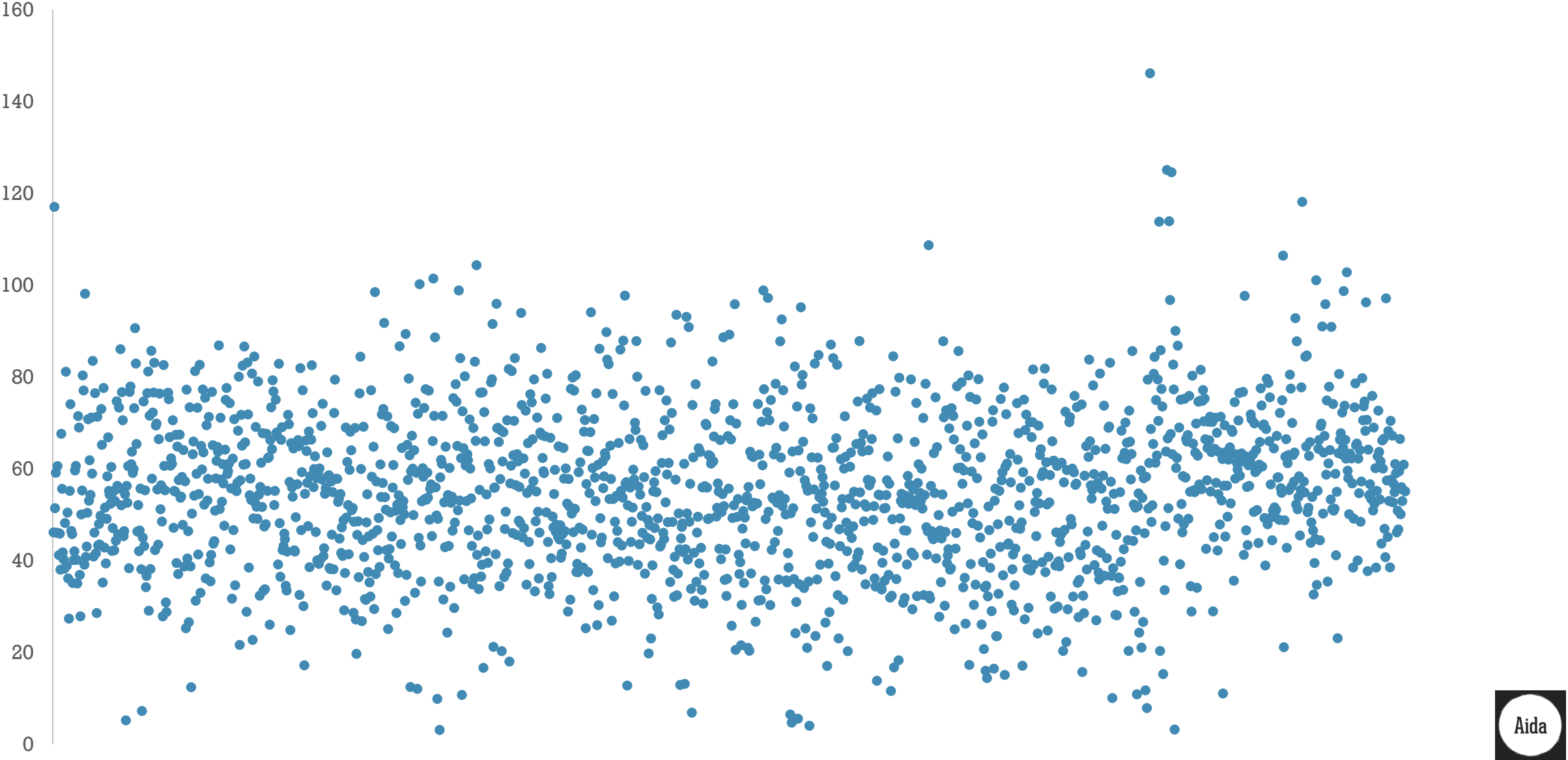
Newspaper page showing
no noisiness

EVALUATION FEATURES

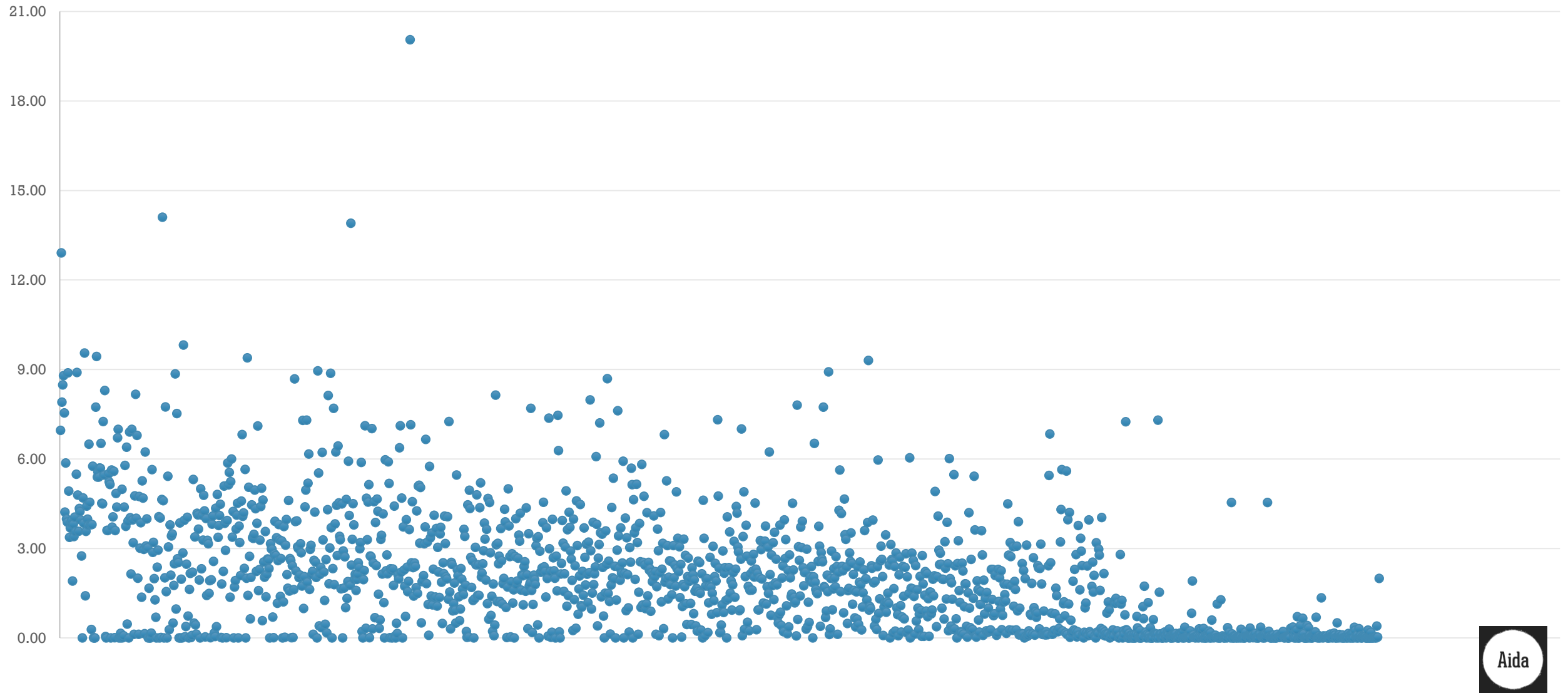
Calculated and analyzed each page and grouped results by language for several features:

- Contrast
- Range effect
- Orientation skew
- Noise

Contrast Test Set



Range Effect Test Set



ANALYSIS:

CONTRAST, RANGE EFFECT

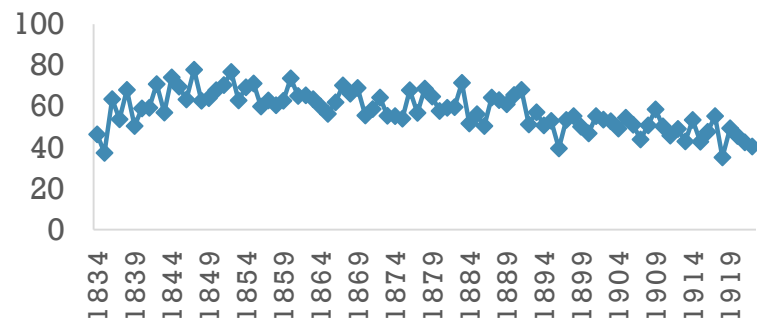
Most newspapers have contrast values between 40 and 80, with contrast > 30 demonstrating good contrast

- The lower the contrast value, the worse the contrast
- $\sim 6-7\%$ images have bad contrast

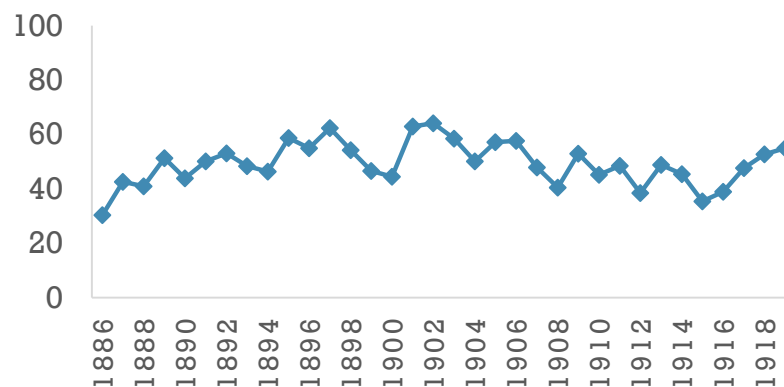
Range effect becomes challenging when > 3

- The ideal value for range effect is 0 (= no range effect is present)
- $\sim 25\%$ images have challenging range effect

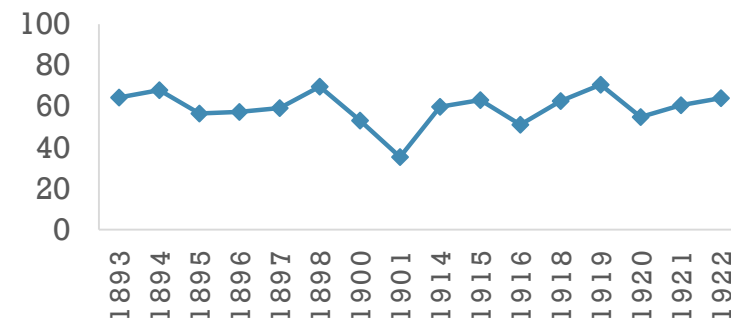
CONTRAST: ENGLISH



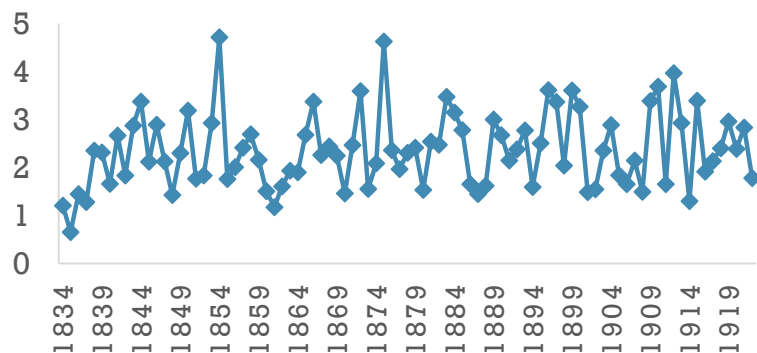
CONTRAST: POLISH



CONTRAST: NORWEGIAN



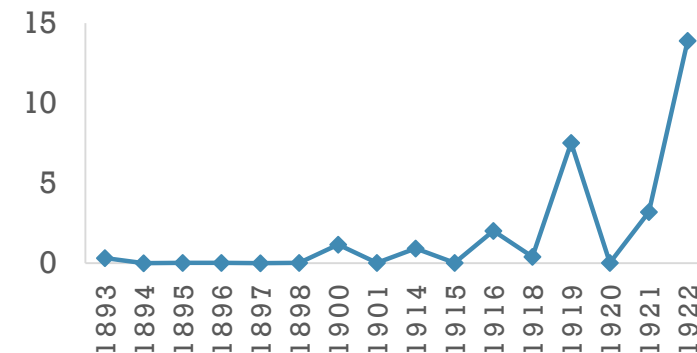
RANGE EFFECT: ENGLISH



RANGE EFFECT: POLISH



RANGE EFFECT: NORWEGIAN



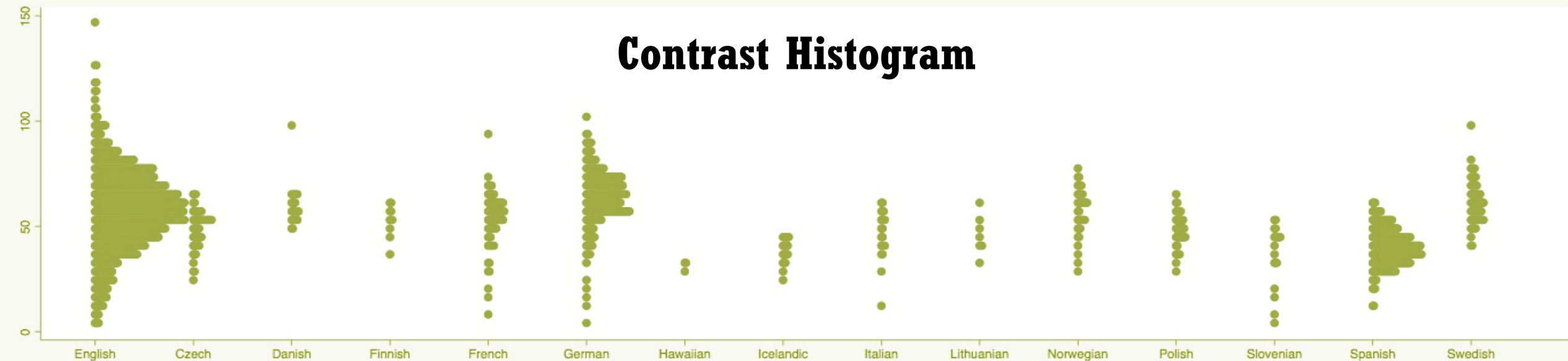
ANALYSIS 2:

CONTRAST, RANGE EFFECT

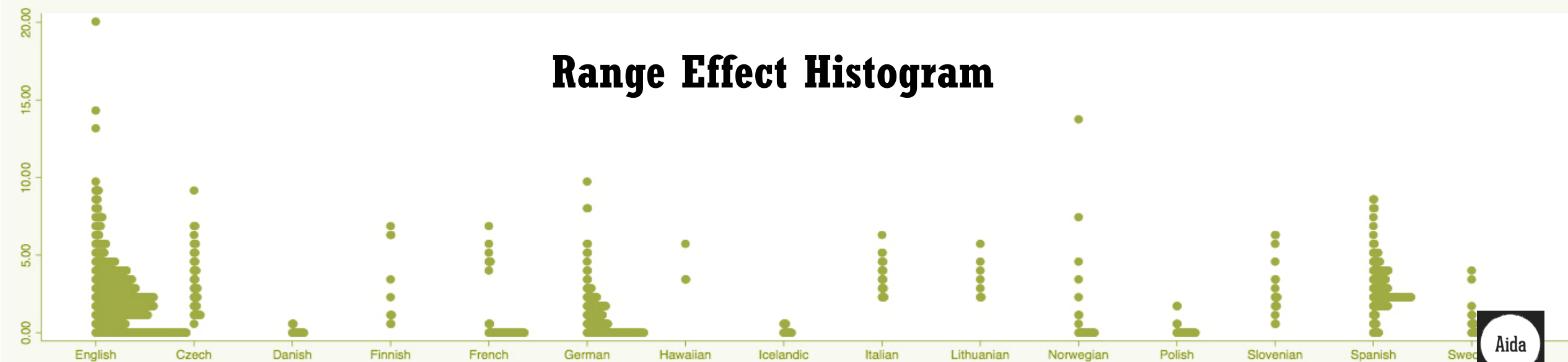
Contrast in all languages is pretty consistent; nor does it change drastically over time

Range effect, on the other hand, not only varies across the different languages, it also changes over time for each language

Contrast Histogram



Range Effect Histogram



ANALYSIS 3:

CONTRAST, RANGE EFFECT

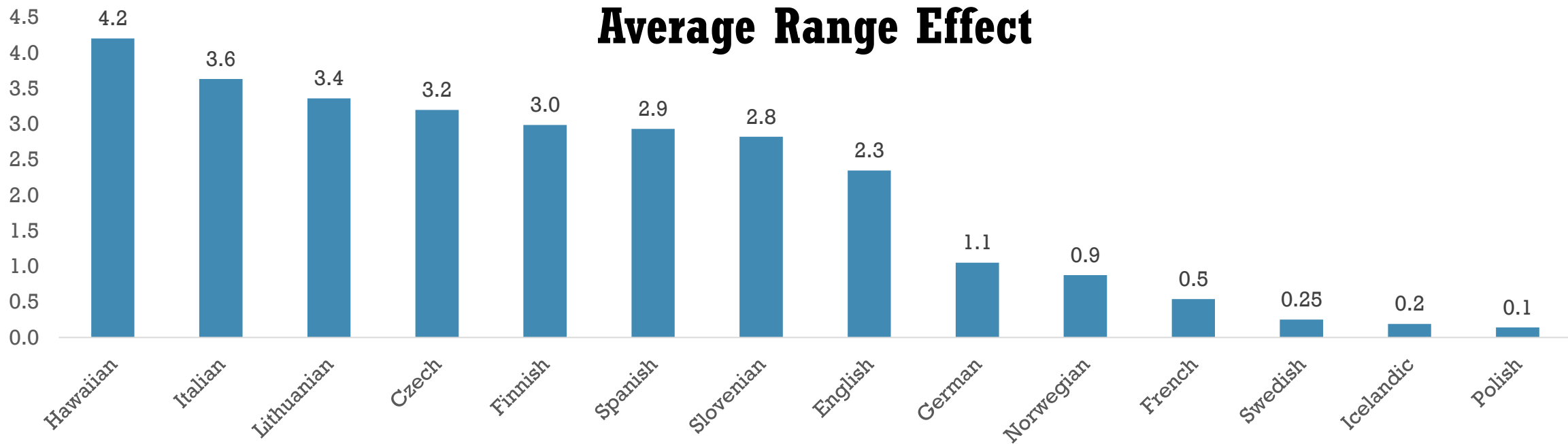
From the histograms:

The data on **Contrast** are fairly **symmetric**, while the **Range Effect** data are **right-skewed**

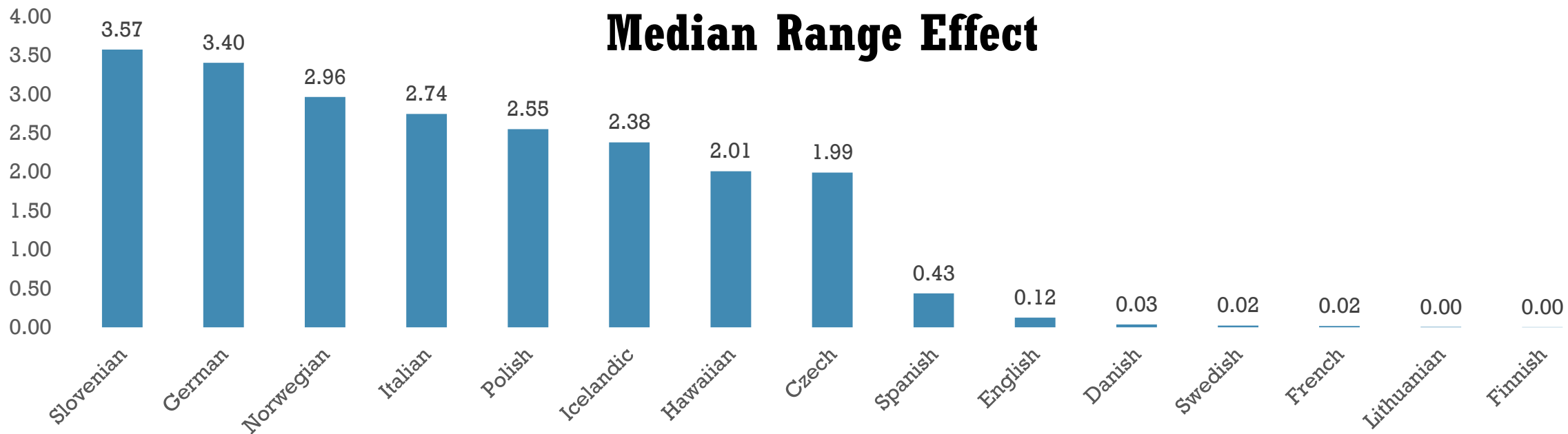
A lot of **outliers** in our Range effect data

- The presence of outliers in our range effect data indicate that the mean (average) would not provide a good estimate: **instead, median would give us a good estimate to evaluate the center of the data**

Average Range Effect



Median Range Effect

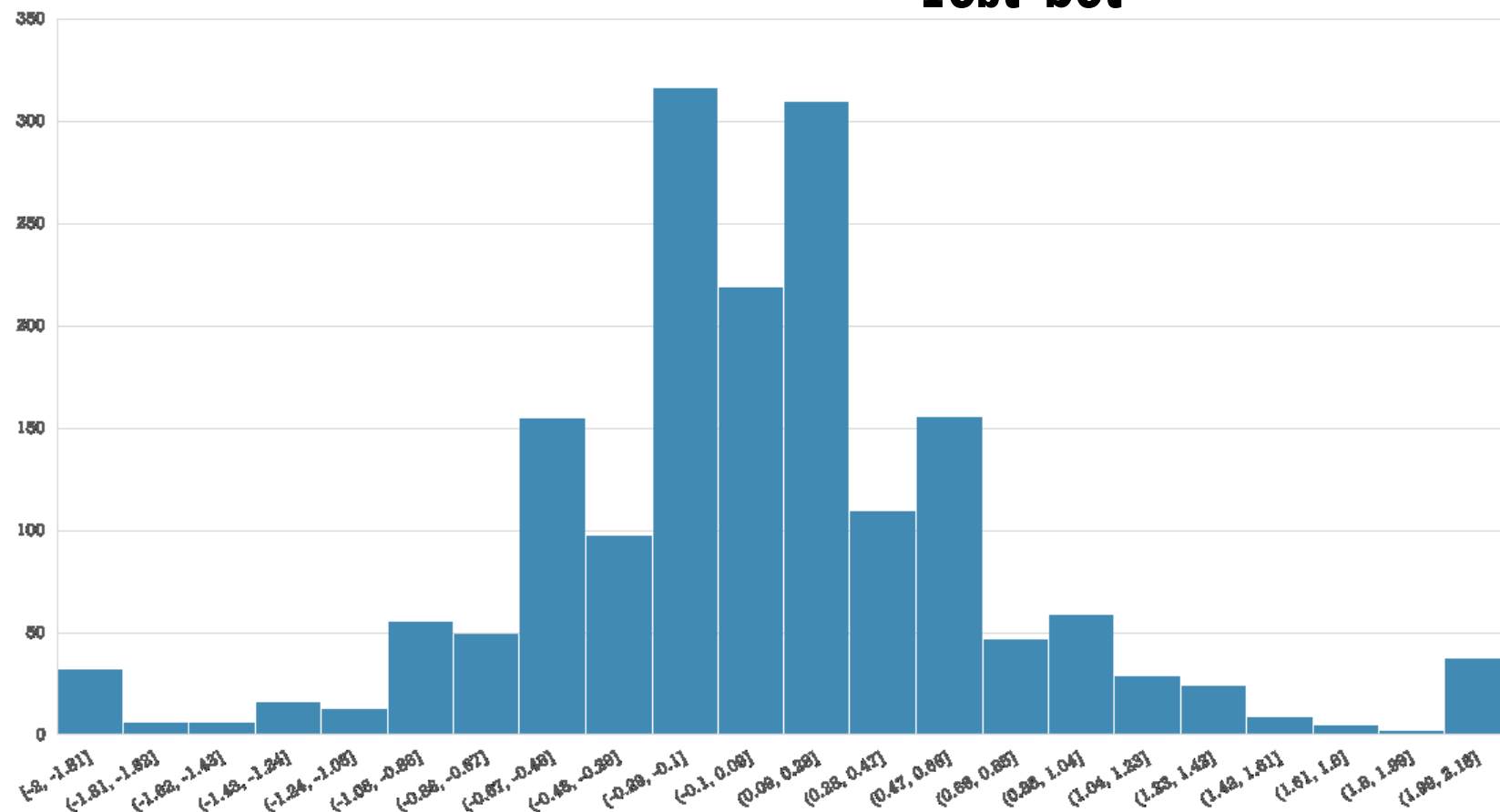


ANALYSIS: ORIENTATION SKEW

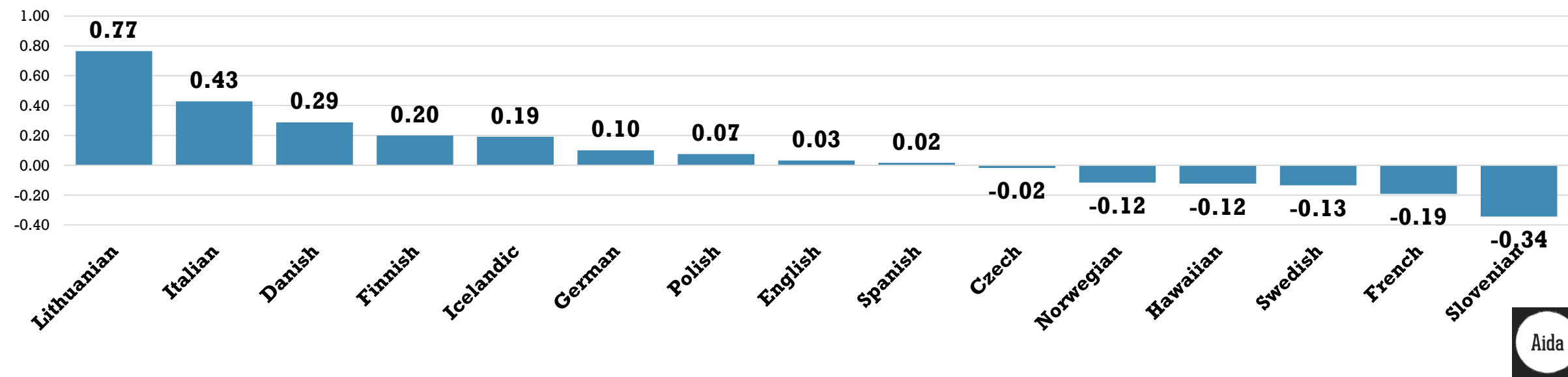
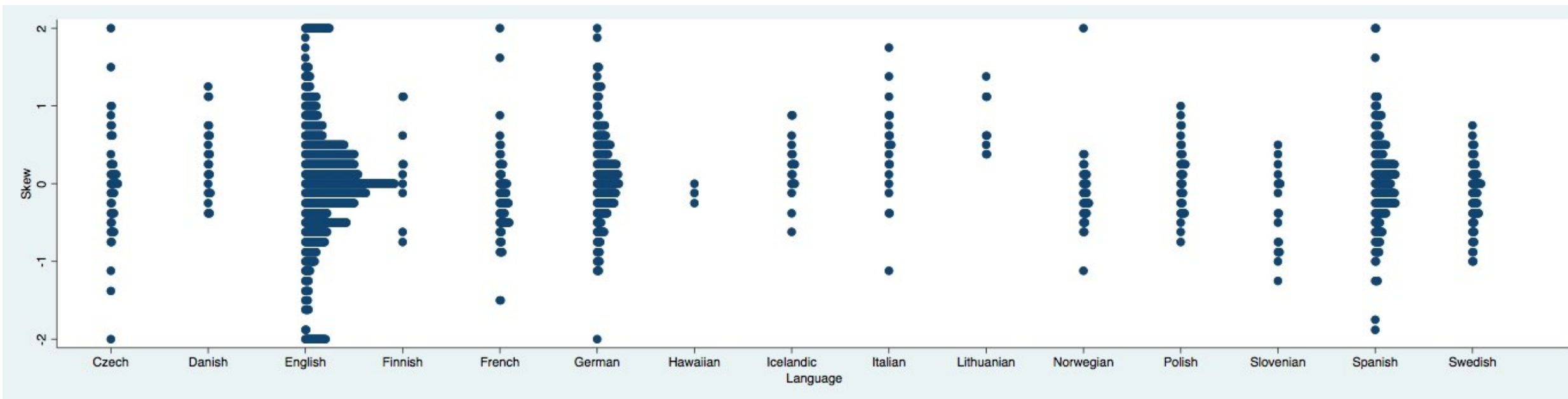
Global orientation skew of the pages in the set

A more effective measure is likely to be local skew, relative to particular parts of the page, or other measures of warpedness or beveled nature of the page.

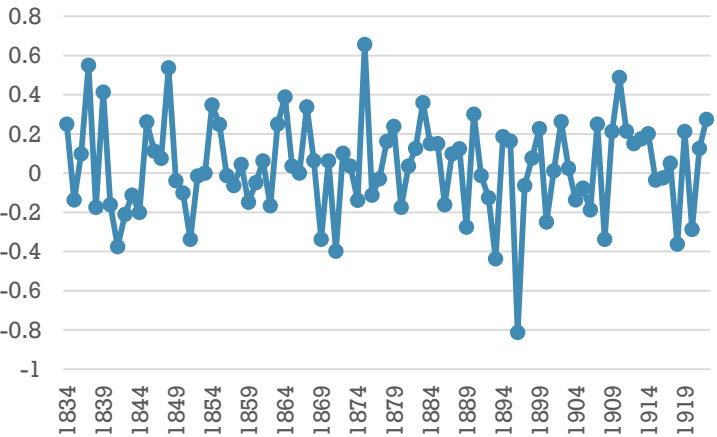
Distribution of Orientation Skew Test Set



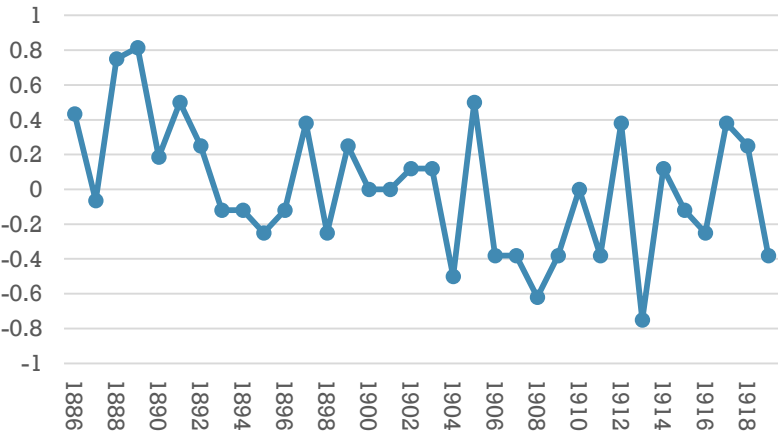
Views of Orientation Skew



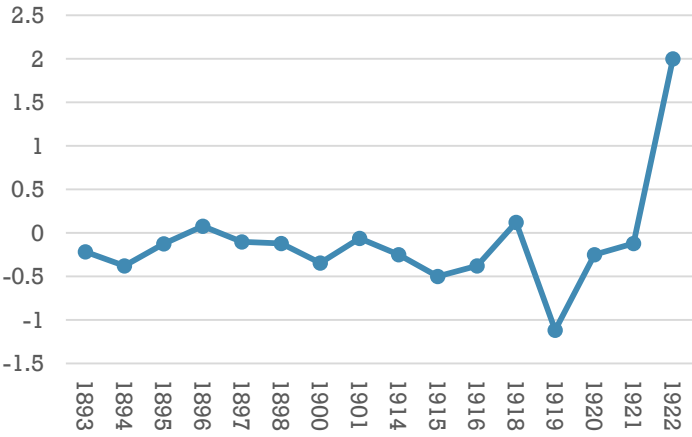
ORIENTATION SKEW: ENGLISH



ORIENTATION SKEW: POLISH



ORIENTATION SKEW: NORWEGIAN





Andrew Prescott

@Ajprescott

Follow



The forgotten labour of digitisation. All newspapers in the British Library were ironed before they were microfilmed to ensure a clear image. Here is a staff member ironing a paper. The microfilm images are used in modern digital packages. [#s482](#) [#KZoo2018](#)



2:01 PM - 12 May 2018

728 Retweets 1,467 Likes



Aida

ANALYSIS: NOISINESS

Assessing effects of bleed through, blobs (e.g., stains), and other non-textual artifacts

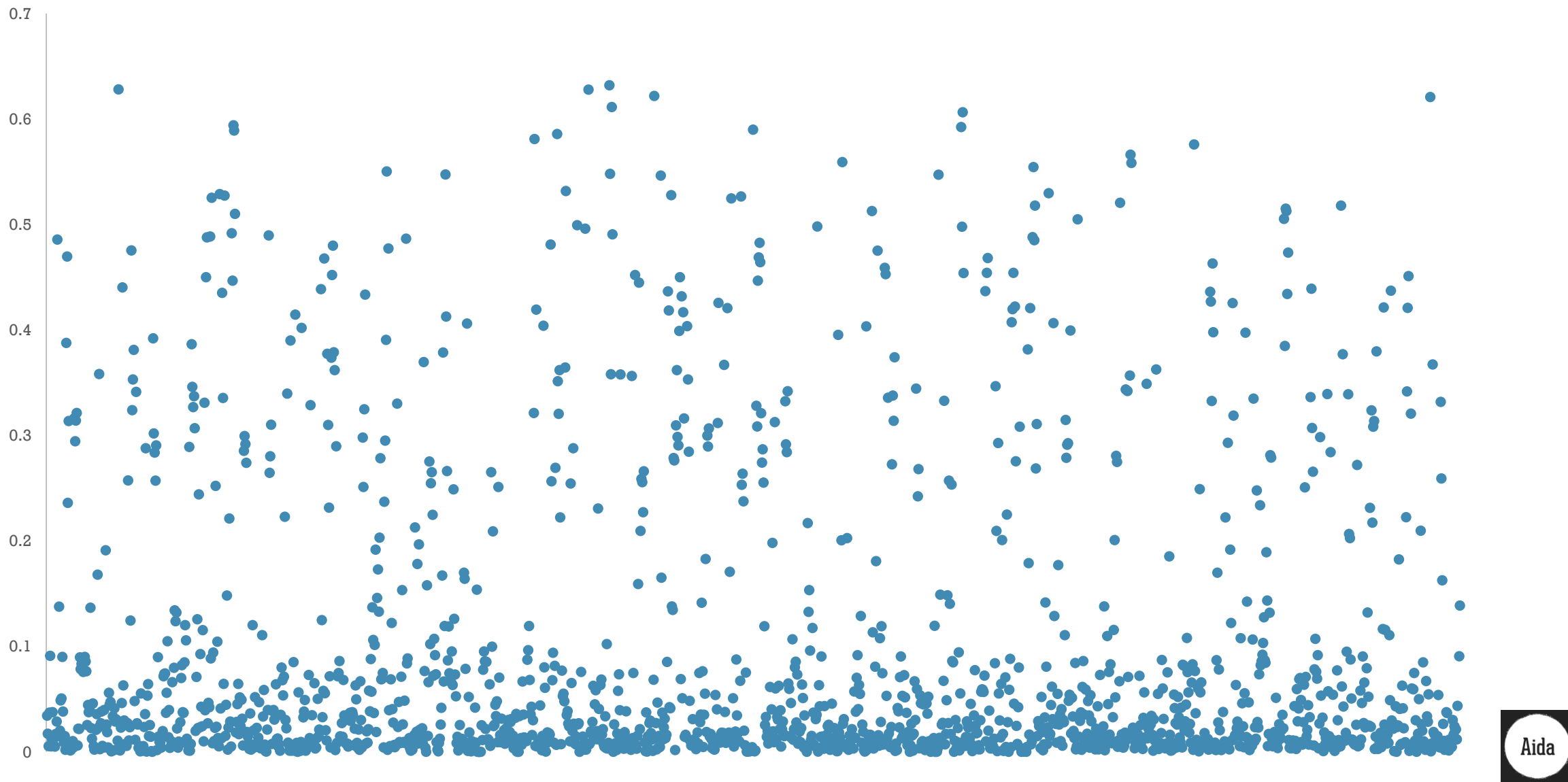
- Defects or degradations of a page, or of the digitization process

Based on histogram analysis—of pixels' intensity values—of each page

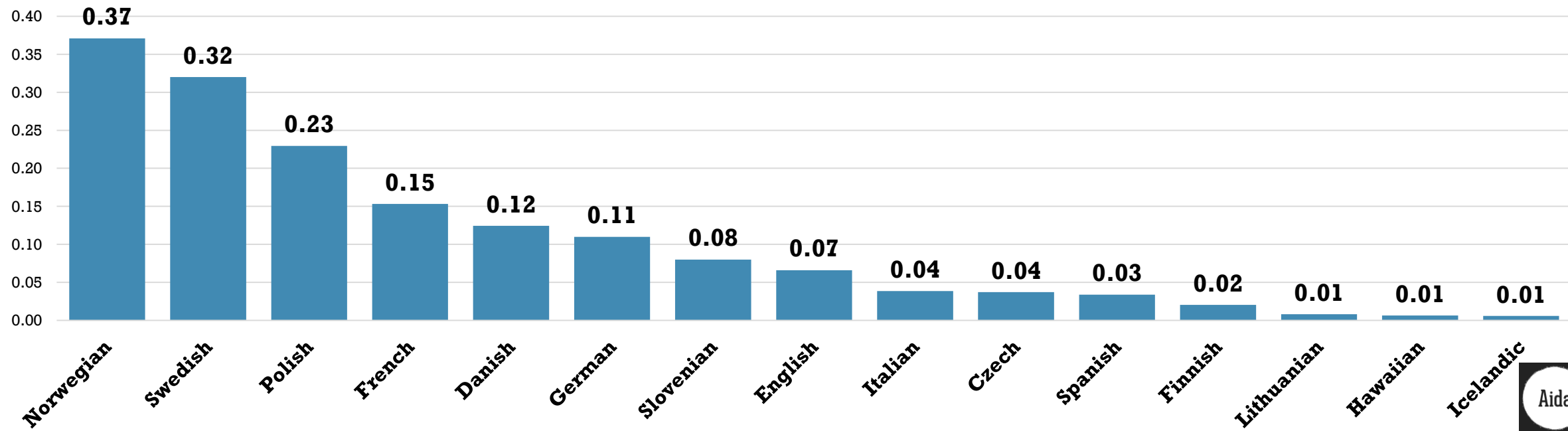
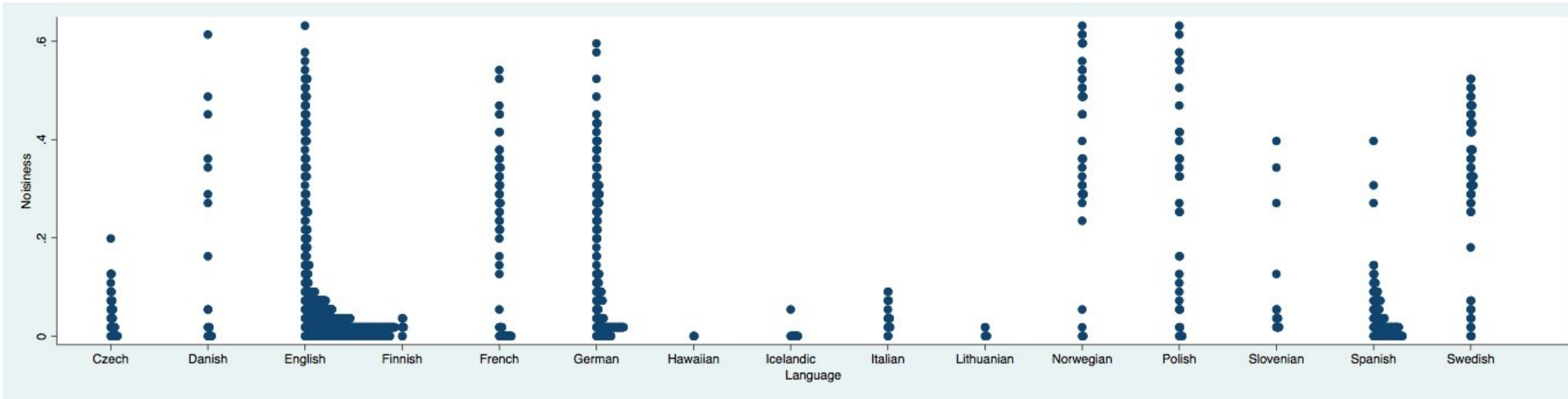
Our current technique:

- If an image is deemed noisy, then it is indeed noisy
- If an image is deemed NOT noisy or of low noise then it is either (1) indeed NOT noisy or of low noise, or (2) of poor contrast

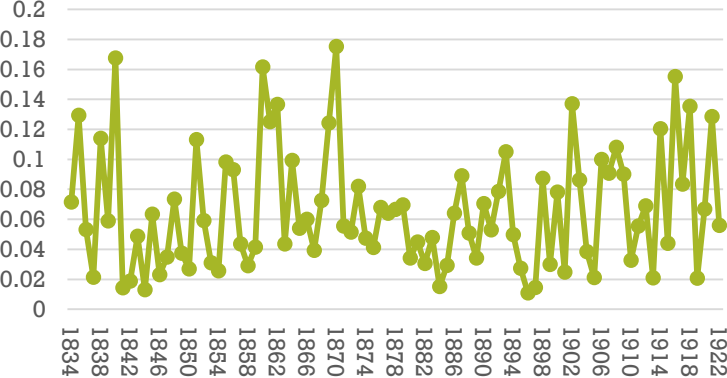
Noisiness Test Set



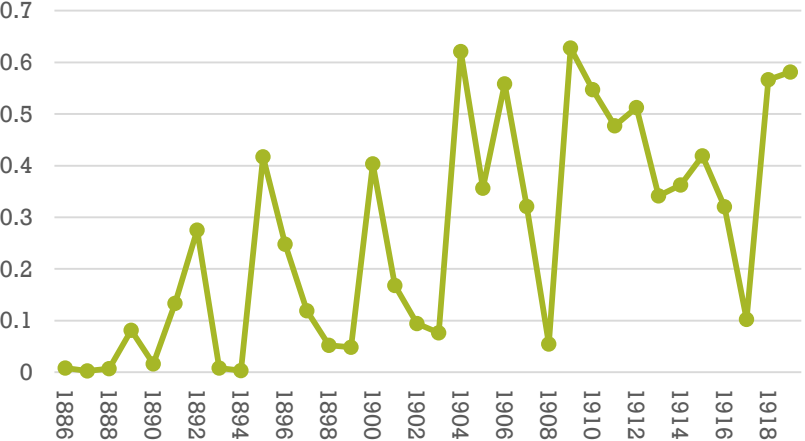
Views of Noisiness



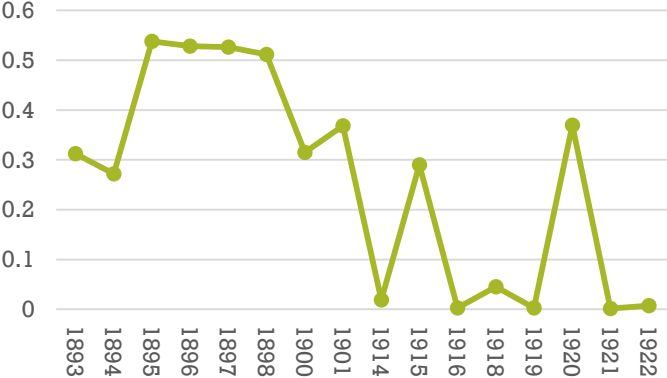
NOISE EFFECT: ENGLISH



NOISE EFFECT: POLISH



NOISE EFFECT: NORWEGIAN

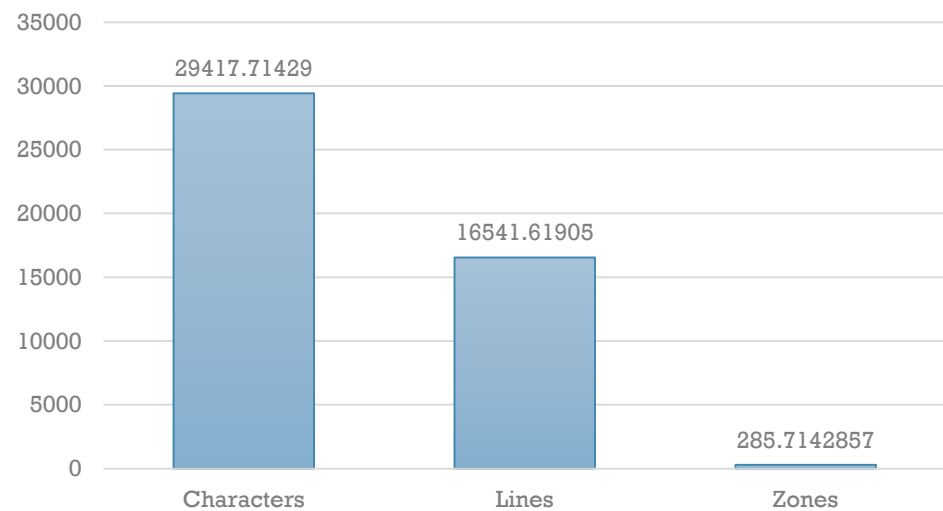


ANALYSIS: CHARACTERS, LINES, ZONES

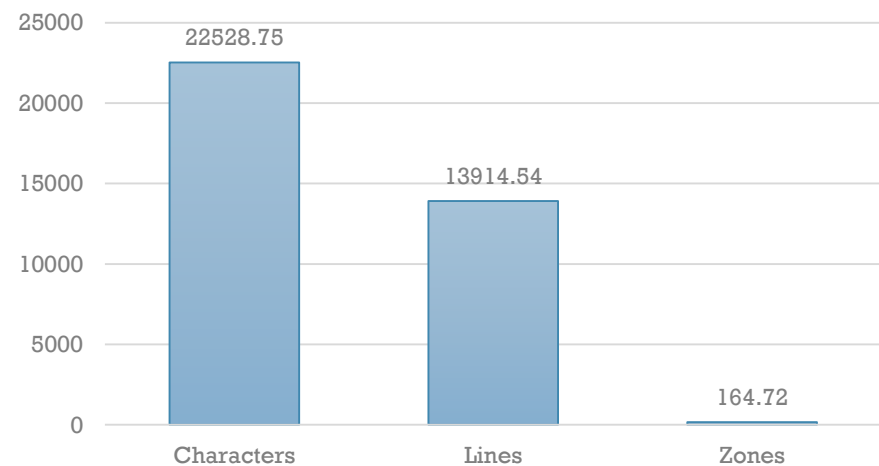
Compute the number of characters, lines, and zones using zoning technique

As a way to begin measuring both complexity and density of the textual content on the page

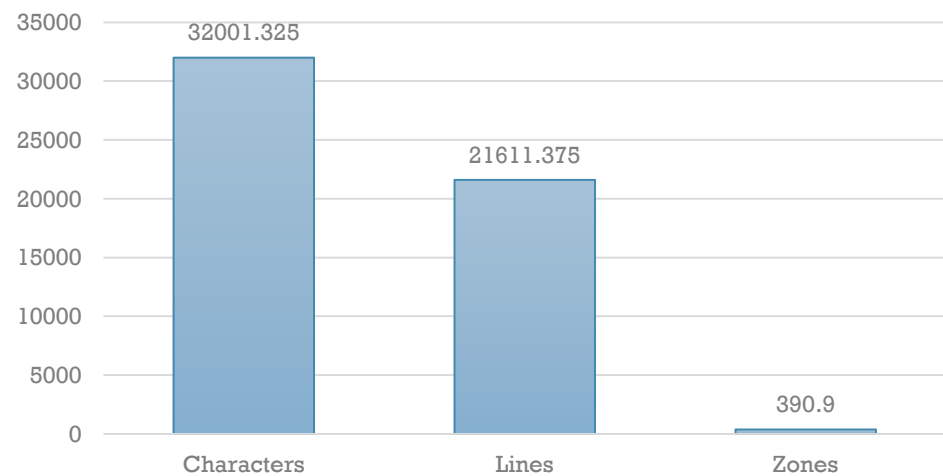
Italian



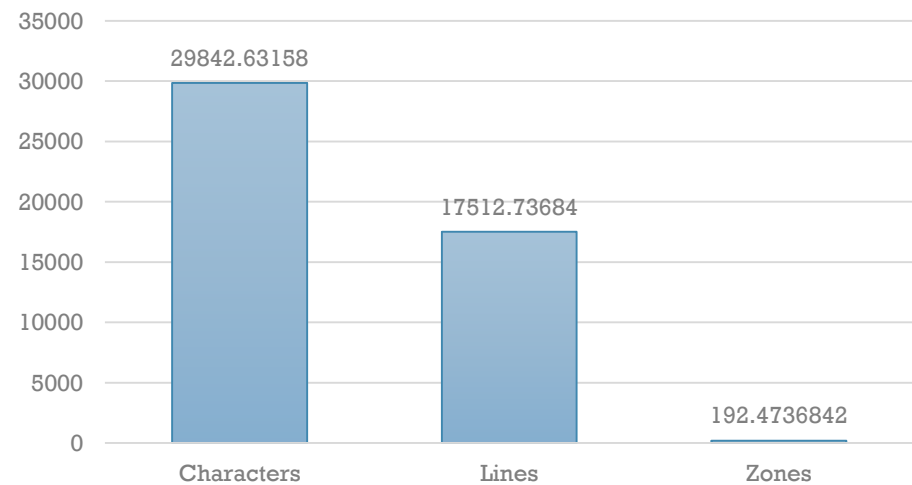
Spanish

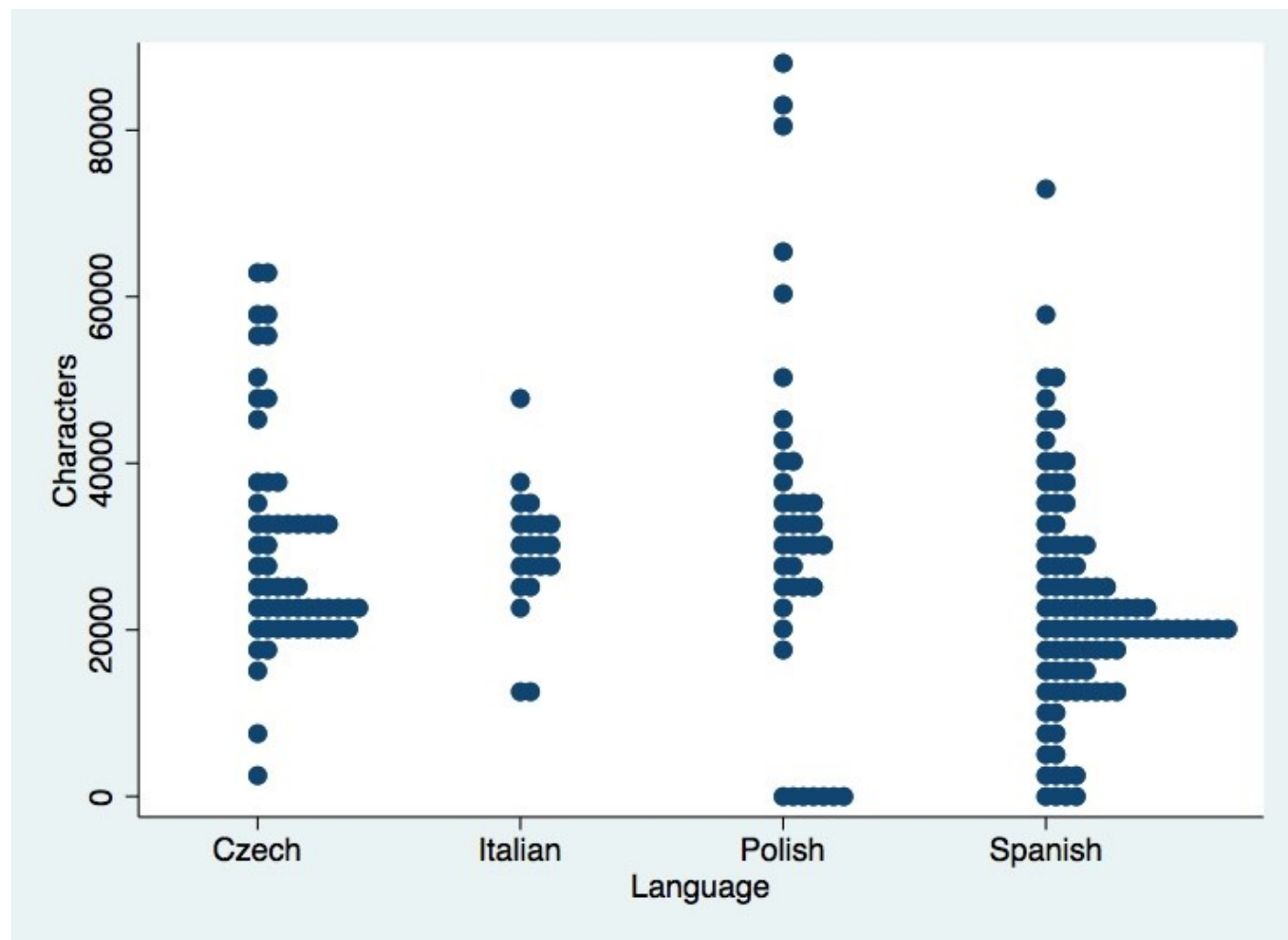


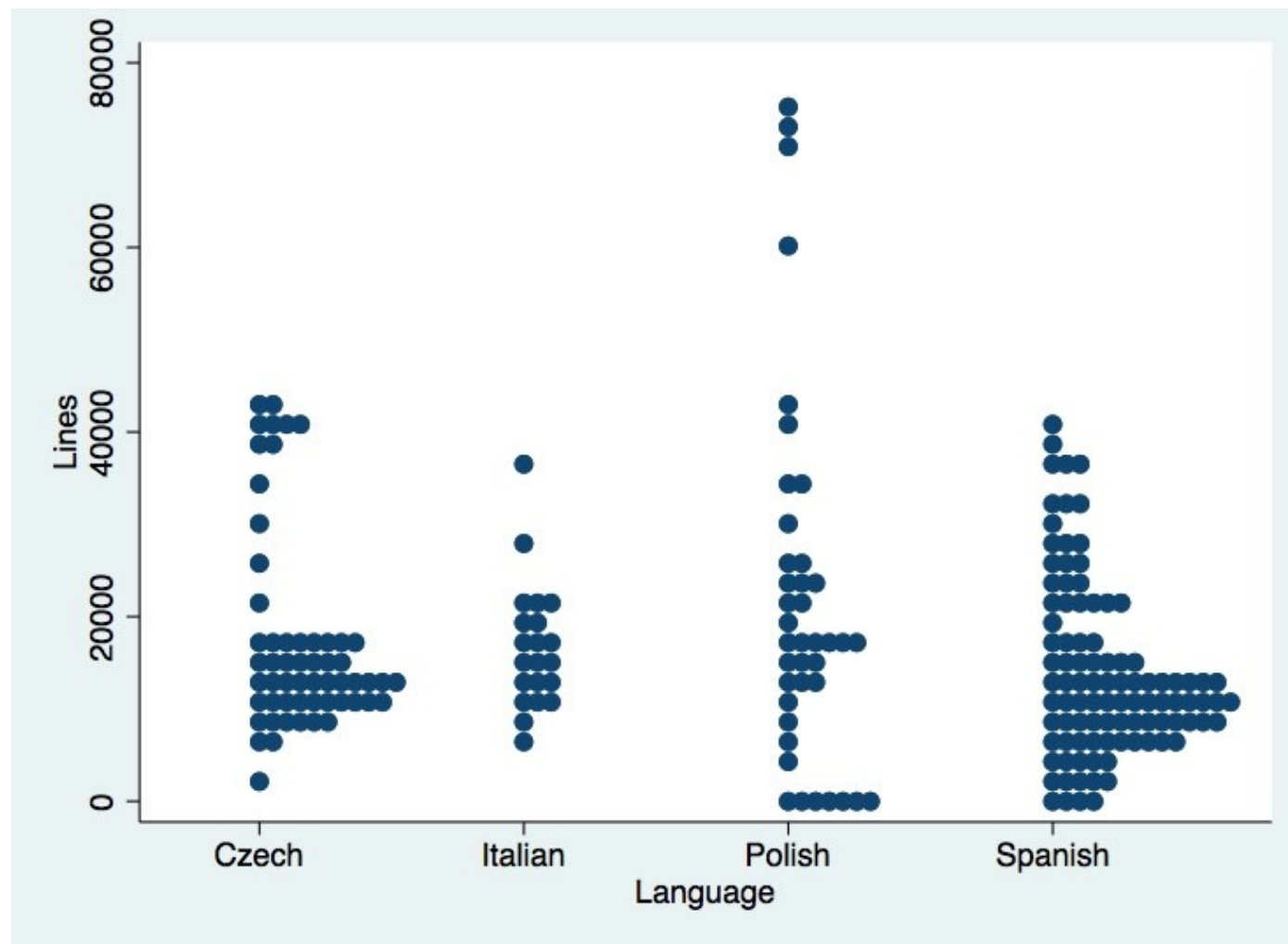
Polish

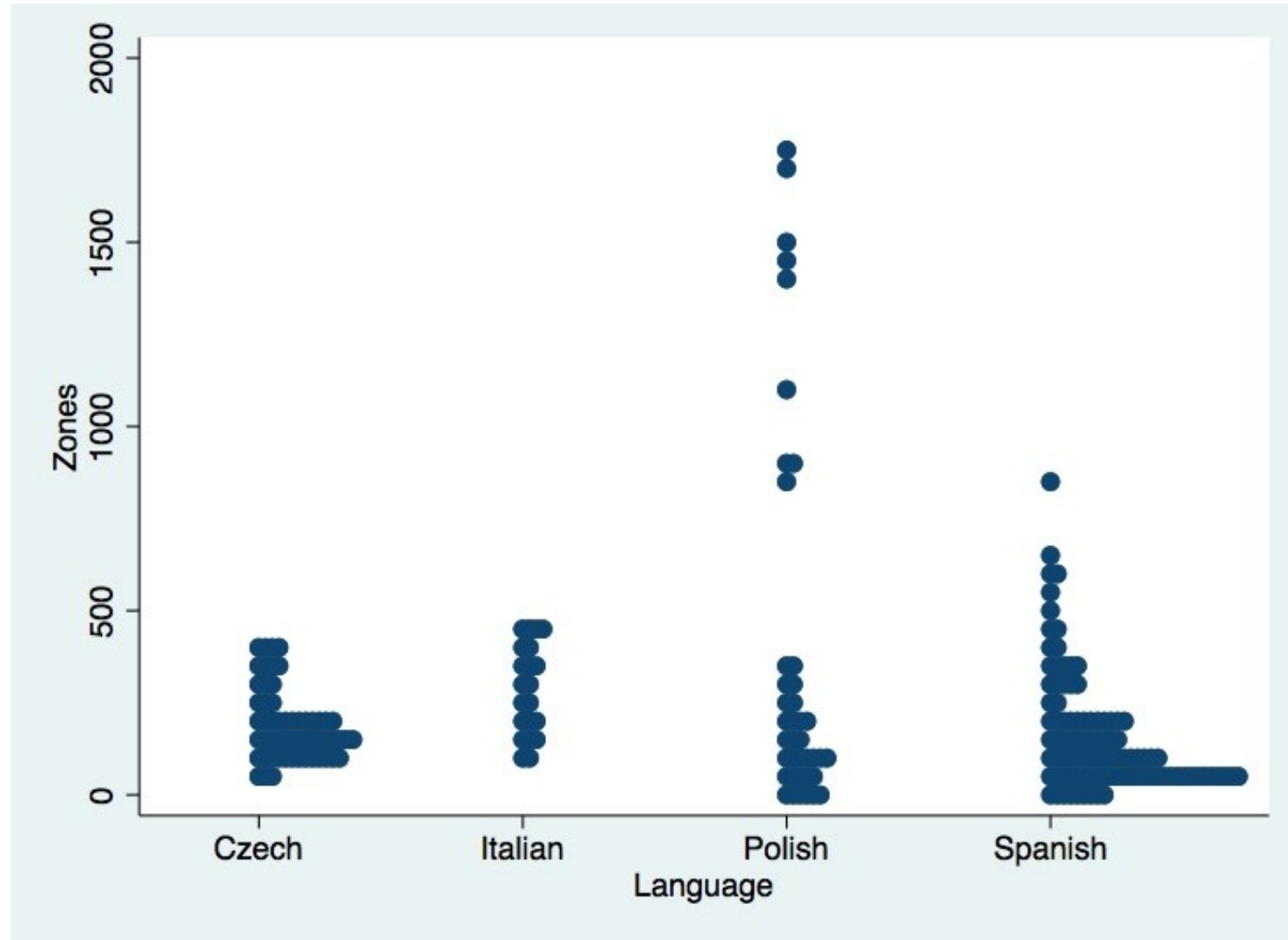


Czech









ANALYSIS 2:

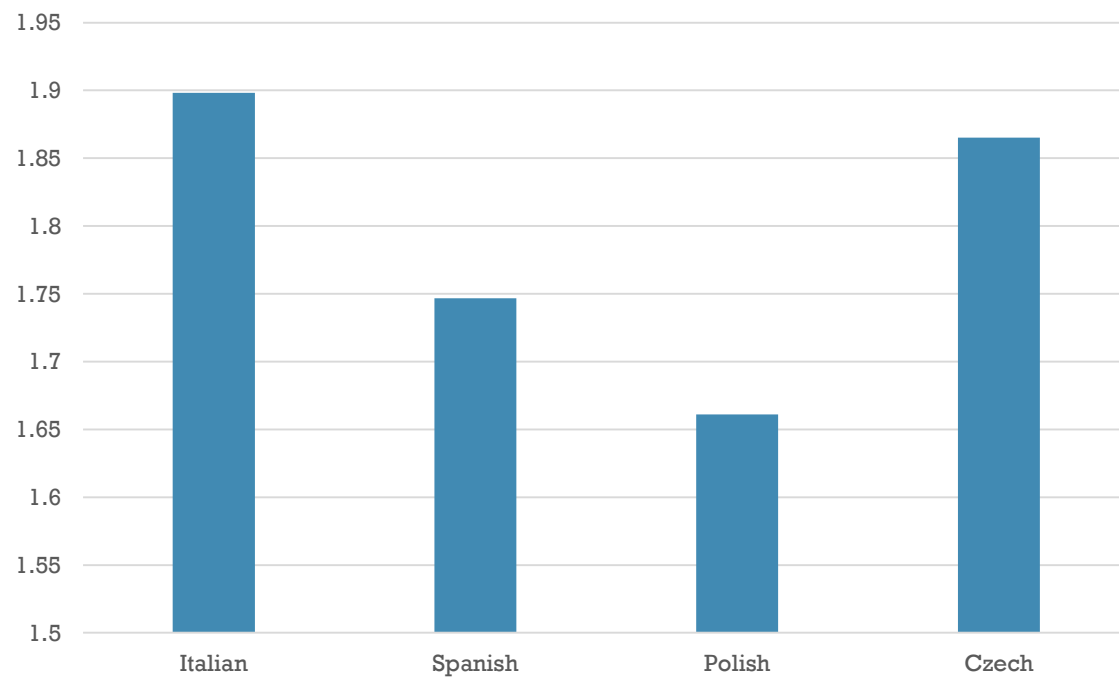
CHARACTERS, LINES, ZONES

Numbers of characters, lines, and zones vary across newspapers of different languages

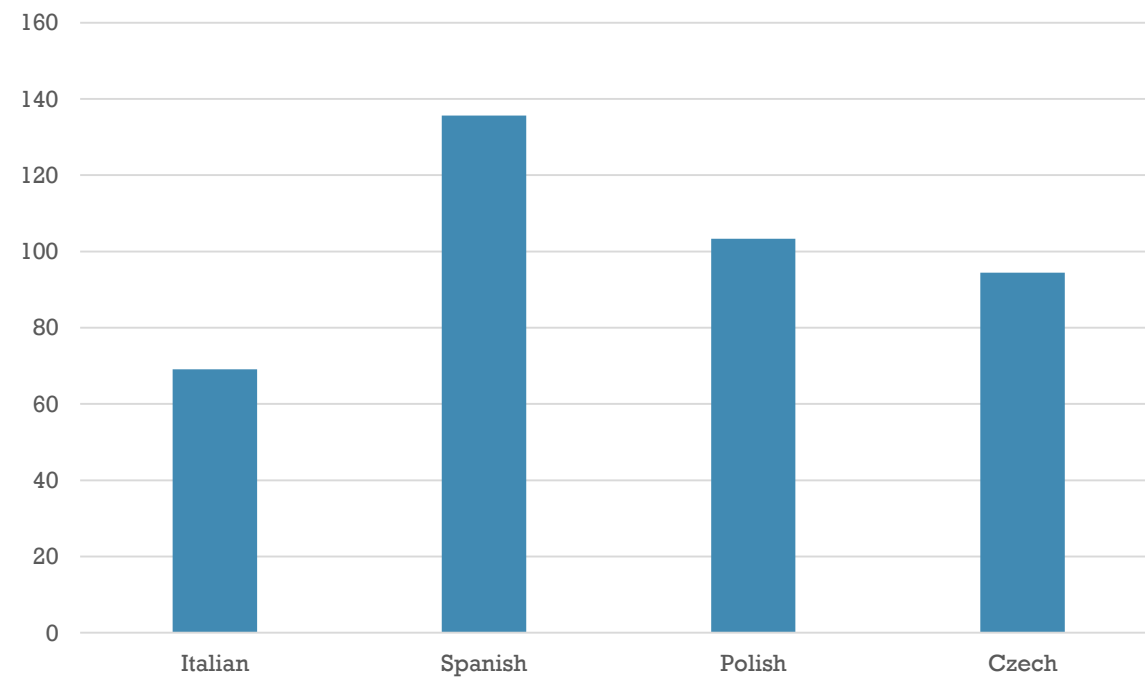
- Polish, Italian, Czech: more characters per page than Spanish
- Italian, Polish: more zones per page than Spanish, Czech

Polish-language newspapers had more outliers; Italian-language newspapers more consistent

Average number of characters per line



Average number of lines per zone



ANALYSIS 3:

CHARACTERS, LINES, ZONES

Italian- and Czech-language newspapers had more characters per line

- Does it mean they had more compact typesetting?
- To investigate further: average length of lines

Spanish-language newspapers had more lines per zone

- Does it mean Spanish-language newspapers had more uniform layouts such that a detected zone was not broken up into multiple smaller zones?
- To investigate further: average size of zones