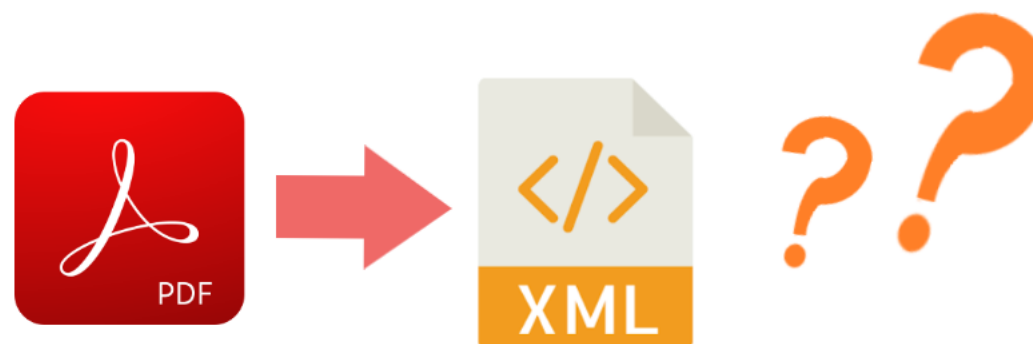# Index Extraction From Textbooks

## KOM Lab Design Workshop – Group KM-2

**Yujin Wang**
guantao0815@gmail.com

**Ruolin Huang**
hrl96fz@foxmail.com

**Fabian Rösch**
fabianroesch@gmail.com

Tim Steuer

# Our Project: Index Extraction From Textbooks

## Motivation

- Index of a book can contain information about its contents that can be analyzed
- If Index can be extracted it can be used for keyword extraction or algorithmic paragraph headline creation

## Problem

- Most Textbooks are PDF
- PDF is not well structured and thus can not easily be interpreted by a machine

## Idea

- It's difficult for a machine to read and contextualize the content of a PDF file
- XML is well structured with labels

## Solution

- Automate index extraction from PDF files to standardized XML files

# Existing Related Projects and Services

**Amazon Textract**
- AWS  cloud based service to extract data from scanned documents using machine learning
- Made for scanned documents that do not contain actual text
- Expensive and not open source

**Apache Tika**
- Java based toolkit to parse PDF Documents and it's metadata and extract into HTML format extraction
- Powerful toolset

**PDFBox**
- Open Source Java Library
- Powerful toolset
- Well documented and integrated into Java

# Implementation Idea: Intermediate Format

## Index

**Symbols**
"Slave Power", 395

**A**
abolitionist, 379, 385
Abu Ghraib, 960
Afghan Northern Alliance, 954
Age of Reason, 114
agrarian society, 20
al-Qaeda, 952, 978
alcalde, 312, 327
Alliance for Progress, 859
American Equal Rights Association, 464
American individualism, 736, 752
American Missionary Association, 457
American Party, 405, 415
American River, 320
American System, 278, 297
Americanization, 499, 505
antebellum, 332, 358
Anti-Federalists, 205, 208
Anti-Imperialist League, 645, 657
Antietam, 429
appeasement, 790
Army of the Potomac, 429, 447
Army of the West, 430, 447
Articles of Confederation, 197
artisan, 269
artisans, 244
Atlanta Compromise, 617, 629

**B**
baby boom, 840, 850

Black Power, 874, 882
Black Pride, 876, 882
black separatism, 875, 882
Black Tuesday, 728, 752
blacklist, 831, 850
Bleeding Kansas, 404, 415
bloody shirt campaign, 577, 598
bonanza farms, 489, 505
Bonus Army, 739, 752
boomerang generation, 975, 978
bootlegging, 712, 719
border ruffians, 402, 415
Boston Harbor, 156
Boston Massacre, 141, 151
Boxer Rebellion, 637
Brains Trust, 758, 783
Bucktails, 276
Bull Run Creek, 426
Bunker Hill, 159
Bush Doctrine, 953, 978

**C**
California Gold Rush, 491, 505
Californios, 317, 327
Calvinism, 44, 59
carpetbagger, 476
carpetbaggers, 470
Carter Doctrine, 914, 915
cash crop, 332, 358
charter schools, 959, 978
chasquis, 13, 30
chattel slavery, 27, 30
checks and balances, 195, 208
chinampas, 12, 30
circuit riders, 363
Citizen Genêt affair, 220, 238
City Beautiful, 560, 567

**Step 1**

```
1   "Slave Power", 395
2   abolitionist, 379, 385
3   Abu Ghraib, 960
4   Afghan Northern Alliance, 954
5   Age of Reason, 114
6   agrarian society, 20
7   al-Qaeda, 952, 978
8   alcalde, 312, 327
9   Alliance for Progress, 859
10  American Equal Rights Association, 464
11  American individualism, 736, 752
12  American Missionary Association, 457
13  American Party, 405, 415
14  American River, 320
15  American System, 278, 297
16  Americanization, 499, 505
17  antebellum, 332, 358
18  Anti-Federalists, 205, 208
19  Anti-Imperialist League, 645, 657
20  Antietam, 429
21  appeasement, 790
22  Army of the Potomac, 429, 447
23  Army of the West, 430, 447
24  Articles of Confederation, 197
25  artisan, 269
26  artisans, 244
27  Atlanta Compromise, 617, 629
28  baby boom, 840, 850
29  bank run, 752
30  bank runs, 730
31  Banking Act of 1935, 774
32  Barnburners, 324, 327
33  Battle of New Orleans, 237
34  Battle of Wounded Knee, 499, 505
35  Bell, 512
36  Beringia, 8, 30
37  bicameral, 203, 208
38  Big Three, 808, 817
39  Bill of Rights, 213, 238
40  Black Cabinet, 778
41  black codes, 458, 476
```

**Step 2**

```xml
<index>
    <file>openSTAX/history/USHistory-OP_tkj0lZo.pdf</file>

    <entries>
    <entry>
        <phrase>"Slave Power"</phrase>
        <pagenumbers>
            <number>395</number>
        </pagenumbers>
    </entry>
    <entry>
        <phrase>abolitionist</phrase>
        <pagenumbers>
            <number>379</number>
            <number>385</number>
        </pagenumbers>
    </entry>
    <entry>
        <phrase>Abu Ghraib</phrase>
        <pagenumbers>
            <number>960</number>
        </pagenumbers>
    </entry>
```

# Evaluation Method: Precision, Recall and F1 Metric

$$\text{Precision} = \frac{Correctly\ Extracted\ Index\ Entries}{Correctly\ Extracted\ Index\ Entries + Falsely\ Extracted\ Index\ Entries}$$

$$\text{Recall} = \frac{Correctly\ Extracted\ Index\ Entries}{Correctly\ Extracted\ Index\ Entries + Falsely\ Considered\ as\ Other\ Text\ Part}$$

$$\text{F1} = 2 * \frac{Precision * Recall}{Precision + Recall}$$

|  |  | Predicted | |
|---|---|---|---|
|  |  | **Negative** | **Positive** |
| **Actual** | **Negative** | True Negative | False Positive |
|  | **Positive** | False Negative | True Positive |

# Status Quo

## Golden Standard

- Five manually created sample XML files for verification and testing

## Programming language and library

- Java
- PDFBox library

## Raw text extraction

- Raw PDF to TXT extraction including all undesired content
- Basic filtering of undesired content and separating phrases from page numbers
- Performs well in one PDF but it still need to be optimized for general using

# Main Obstacles Right Now

- Determining where the Index in any given PDF is

- Filtering and removal of non-relevant content (e.g. page number, header, etc)

- How to identify subphrases



Headlines

These are not individual entries but subentries

Header /Footer

Pagenumber

# Further Problems: Special Cases

- Some special symbols
- The phrase is too long and split into two lines
- How to distinguish between the hyphen existing in the phrase itself and the hyphen due to line breaks

Line break due to long entrys

special symbols

hyphens due to line breaks

# Sources

[1]  Corbett, P. S., Volker, J., Lund, J. M., Pfannestiel, T., Vickery, P. S., & Janssen, V. (2014). *U. S. History*. Amsterdam, Niederlande: Amsterdam University Press.

[2]  Koo Ping Shung, (Mar 15, 2018) Accuracy, Precision, Recall or F1? https://towardsdatascience.com/accuracy-precision-recall-or-f1-331fb37c5cb9

[3]  Kozen, D. C. (2013). *Automata and Computability*. New York, Vereinigte Staaten: Springer Publishing.

[4]  Flowers, P., Theopold, K., Langley, R., STEPHEN F., ROBINSON, W. R. (2015). Chemistry-OP. Rice University Press

[5]  RYE, C., WISE, R., JURUKOVSKI, V., DESAIX, J., CHOI, J., AVISSAR, Y. (2013). Biology-OP. Rice University Press

[6] Olivier Bonaventure (October 30, 2011) *Computer-Networking-Principles-Bonaventure-1-30-31-OTC1.*