

# Digitizing Historical Balance Sheet Data: A Practitioner's Guide

---

Sergio Correia<sup>a</sup>    Stephan Luck<sup>b</sup>

December 3, 2021

Methodological Advances in the Extraction and Analysis of Historical Data

---

<sup>a</sup>Federal Reserve Board, [sergio.a.correia@frb.gov](mailto:sergio.a.correia@frb.gov)

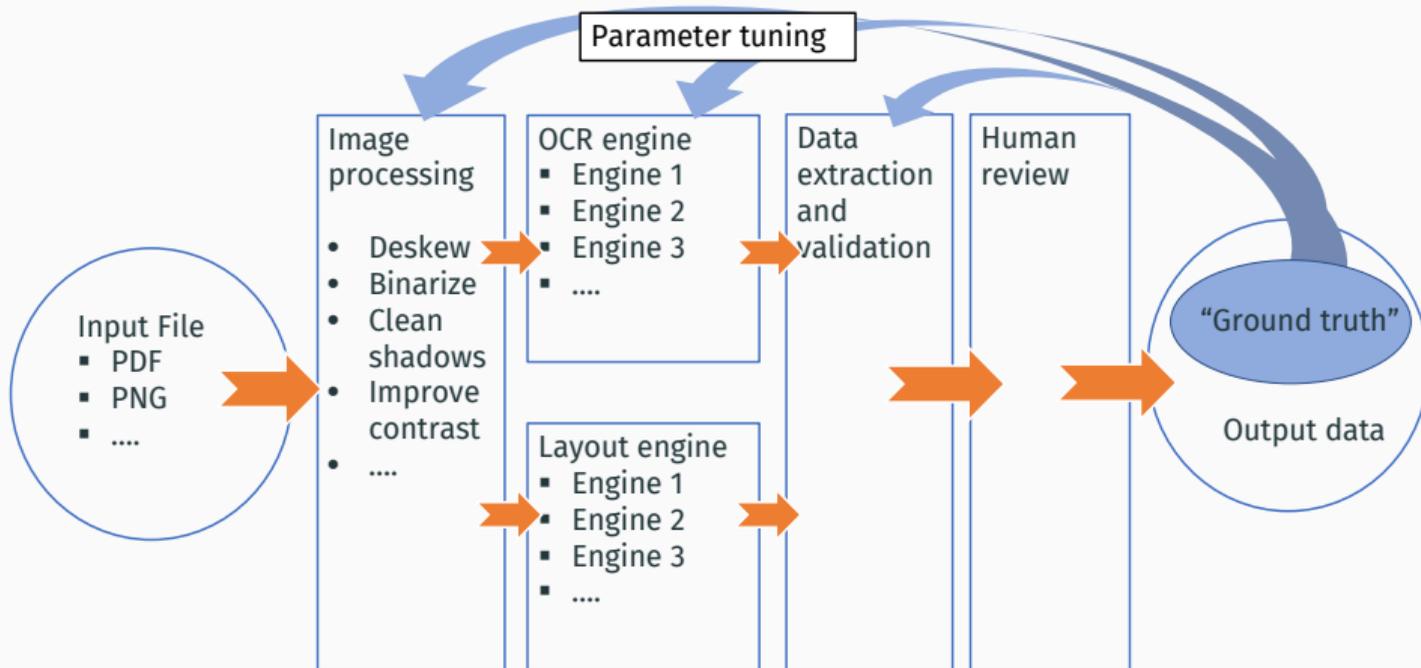
<sup>b</sup>Federal Reserve Bank of New York, [stephan.luck@ny.frb.org](mailto:stephan.luck@ny.frb.org)

## Goals of this talk (and paper)

- Share what we learned from transcribing structured historical data,
  - *At scale*
  - With limited resources (no funding, no RAs or interns)
- Show,
  - Our mistakes so you can avoid them
  - What we think works well (beyond just our projects)
- Make digitization projects more accessible:
  - Key principle: modular tools that can be mix-and-matched as needed
  - You don't have to reinvent the wheel, you *don't* have to be an expert coding
  - Open sourcing our tools: <https://github.com/sergiocorreia/quipucamayoc>

- Off-the-shelf OCR tools,
  - Have improved dramatically in recent years
  - Potentially very powerful
  - But still have high error rates, particularly for historical data
  - This makes large-scale efforts unfeasible
- Insight: we need to augment OCR with pre- and post-processing methods
- *The researcher must become the practitioner*
  - Understand their data
  - How to work around its limitations
  - How to take advantage of its characteristics

# The data extraction pipeline



## Example #1: OCC Annual Reports to Congress

- More than 100,000 balance sheets in *tabular form* (1867-1904)
- About 12,000 hours of work if typed by hand
- Used in Carlson et al. JPE, Forthcoming

**NEW YORK.**

**First National Bank, Watertown.**

EDWIN L. PADDOCK, *President.*                      No. 73.                      OSCAR PADDOCK, *Cashier.*

Resources.		Liabilities.	
Loans and discounts .....	\$21,402 33	Capital stock paid in .....	\$100,000 00
Overdrafts .....		Surplus fund .....	20,000 00
U. S. bonds to secure circulation .....	100,000 00	Other undivided profits .....	10,325 45
U. S. bonds to secure deposits .....		National bank notes outstanding ..	86,915 00
U. S. bonds on hand .....	45,000 00	State bank notes outstanding .....	
Other stocks, bonds, and mortgages ..	17,200 00	Dividends unpaid .....	
Due from approved reserve agents ..	14,405 83	Individual deposits .....	10,811 99
Due from other banks and bankers ..		United States deposits .....	
Real estate, furniture, and fixtures ..	21,000 00	Deposits of U. S. disbursing officers ..	
Current expenses and taxes paid ..		Due to other national banks .....	
Premiums paid .....		Due to State banks and bankers .....	
Checks and other cash items .....	1,283 04	Notes and bills re-discounted .....	
Exchanges for clearing-house .....		Bills payable .....	
Bills of other banks .....	752 00		
Fractional currency .....			
Specie .....	2,509 24		
Legal-tender notes .....			
U. S. certificates of deposit .....			
Due from U. S. Treasurer .....	4,500 00		
<b>Total .....</b>	<b>228,052 44</b>	<b>Total .....</b>	<b>228,052 44</b>

## Example #2: Saling's Börsen-Papiere

- More than 30,000 balance sheets and income statements of German financial and non-financial firms (1915-1933)
- Used in ongoing research with Brunnermeier and Zimmermann



## Example #2: Saling's Börsen-Papiere

- Data is in **FREE FORM TEXT** !!!
  - No predefined set of labels
  - Archaic abbreviations (in German)
  - Values twelve-digit long in hyperinflation years

bank, Darmstädter u. Nationalbank, Delbrück Schickler & Co., Disconto-Ges., Dresdner Bank, Hardy & Co., Hugo J. Herzfeld; Leipzig: Allgem. Deutsche Credit-Anstalt; Frankfurt a. M.: Metallbank u. Metallurgische Ges. — *Kurs*: Freihänd. Verkauf zu 98%. Zugel. Jan. 1923. — Erster Kurs 8./1. 1923: 99%. — Ult. 1923: 200%. (Auch in Leipzig notiert).

*Gewinn* 1922: Zuschuss d. Mansfeld A.-G. 318 200 213 *ℳ*. — *Dagegen*: Betriebsverlust 286 714 717, Verlust a. Kolonie 8 339 701, Zinsen 9 925 385, allgem. Unk. 22 379 195, Anleihe- do. 9 539 564, Abschreib. 1 301 650 *ℳ*.

*Bilanz* ult. 1922: *Aktiva*: Gerechsamte 860 000, Grundstücke 370 636, Schächte 3 224 658, Betriebsgebäude 3 368 342, Maschinen- u. Dampfkessel 3 119 166, Wasserhaltungs-Anlagen 25 377, elektr. Licht- u. Kraftanlagen 327 568, Bahnanlagen 624 321, Wege, Zechenplatz, Be- u. Entwässerung 107 996, Wohn- u. Wirtschaftsgeb. 63 914, Betriebsgeräte 151 038, Mobilar 1399, Kokerei-Betriebsgebäude 69 010 (zus. 12 313 423), Kasse 397 416, Magazinbestand 20 431 136, Kohlenbestand 758 690, Beteilig. 133 750, Aussenstände: Guth. bei d. Mansfeld A.-G. 75 184 127, Vorauszahlungen 118 787 325, Sonstige 506 130 073 (zus. 698 081 526); Verlust 318 200 213 *ℳ*. — *Passiva*: Kapital 2 000 000, Anleihen: Ausgabe vom Jahre 1914 I. Einzahl. 5% 4 692 500, II. Einzahl. 4½% 10 000 000, Ausgabe vom Jahre 1920 5% 40 000 000, do. v. Jahre 1922 5% 100 000 000 (zus. 154 692 500); Anl.-Zinsen 2 459 784, Hyp. 8500, Akzepte 165 000 000, Schulden: Bankschulden 5 198 753, Guthaben d. Mansfeld A.-G. 8 324 700, Sonstige 712 631 828 (zus. 726 155 281) *ℳ*. — S. b. 1 Md. *ℳ*.

### **Sächsische Gussstahl-Werke Döhlen Akt.-Ges. in Dresden.**

(Bis 27./10. 1920: *Sächsische Gussstahlfabrik in Döhlen*).

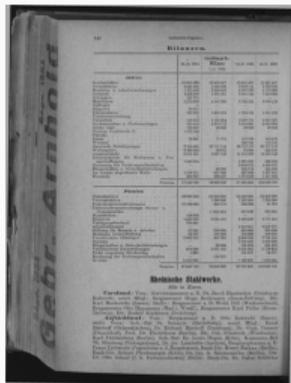
*Vorstand*: Gen.-Dir. Kommerz.-R. Herm. Pfeifer, Gen.-Dir. Kommerz.-R. Herm. Boehm. — *A.-R.*: Vors.: Geh. Kommerz.-R. Gen.-Konsul Gust. v. Klempener; Stellv.: Kommerz.-R. Willy Osswald [Deutsche Bk.]; sonst. Mitgl.: Bergtrat a. D. Andreas Nägel, Justizrat Dr. Jul. Bondi, Geh. Baurat Fritz Martiny (Ob.-Schreiberhau), Bankier Dr. Gustaf Ratjen (Berlin-Dahlem), Kfm. Albert Roth-

## Tackling the datasets

- Each dataset required a different combination of tools
- OCC: good quality scan, already black-and-white
  - No need for preprocessing
  - But required lots of layout recognition for the tables
- *Saling's*: scan not as good, in color, often with artifacts
  - Lots of preprocessing
  - Less need for layout recognition
- Also, what are the *advantages* of these datasets?
  - Balance sheets allow us to validate data: sum of assets = total assets; assets=liabilities



# Showcase - Fixing Size Distortions



(a) Input



(b) Remove noise



(c) Dilate white area



(d) Identify white box



(e) Trim





# Showcase - Human Review Interfaces

**NEW JERSEY National Union Bank, Down.**

Assets	Liabilities
Loans and discounts: \$95,141.85	Capital stock paid in: \$95,000.00
Real estate: 2,000.00	Surplus fund: 30,000.00
U.S. bonds to secure deposits: 20,000.00	Other undivided profits: 12,000.00
U.S. bonds on hand: 10,000.00	Reserve for cash on hand: 41,000.00
Other assets, bank, and mortgage: 10,000.00	From bank notes outstanding: 407.00
Due from approved reserve agents: 40,000.00	Overdrafts unpaid: 407.00
Due from other banks and bankers: 20,000.00	Individual deposits: 47,000.00
Due from other banks and bankers: 20,000.00	Deposits for U.S. Savings certificates: 47,000.00
Prepaid rent: 1,000.00	Due to other national banks: 5,000.00
Prepaid taxes: 1,000.00	Due to other banks and bankers: 1,000.00
Prepaid interest: 1,000.00	Notes and bills to be advanced: 1,000.00
Prepaid interest: 1,000.00	Other payables: 1,000.00
Prepaid interest: 1,000.00	Total: 200,000.00

**First National Bank, Elizabeth.**

Assets	Liabilities
Loans and discounts: \$200,000.00	Capital stock paid in: \$200,000.00
Real estate: 10,000.00	Surplus fund: 40,000.00
U.S. bonds to secure deposits: 10,000.00	Other undivided profits: 30,000.00
U.S. bonds on hand: 10,000.00	Reserve for cash on hand: 71,000.00
Other assets, bank, and mortgage: 10,000.00	From bank notes outstanding: 300.00
Due from approved reserve agents: 10,000.00	Overdrafts unpaid: 300.00
Due from other banks and bankers: 10,000.00	Individual deposits: 600,000.00
Due from other banks and bankers: 10,000.00	Deposits for U.S. Savings certificates: 600,000.00
Prepaid rent: 1,000.00	Due to other national banks: 5,000.00
Prepaid taxes: 1,000.00	Due to other banks and bankers: 1,000.00
Prepaid interest: 1,000.00	Notes and bills to be advanced: 1,000.00
Prepaid interest: 1,000.00	Other payables: 1,000.00
Prepaid interest: 1,000.00	Total: 1,000,000.00

**National State Bank, Elizabeth.**

Assets	Liabilities
Loans and discounts: \$500,000.00	Capital stock paid in: \$500,000.00
Real estate: 10,000.00	Surplus fund: 100,000.00
U.S. bonds to secure deposits: 10,000.00	Other undivided profits: 100,000.00
U.S. bonds on hand: 10,000.00	Reserve for cash on hand: 300,000.00
Other assets, bank, and mortgage: 10,000.00	From bank notes outstanding: 100,000.00
Due from approved reserve agents: 10,000.00	Overdrafts unpaid: 100,000.00
Due from other banks and bankers: 10,000.00	Individual deposits: 1,000,000.00
Due from other banks and bankers: 10,000.00	Deposits for U.S. Savings certificates: 1,000,000.00
Prepaid rent: 1,000.00	Due to other national banks: 5,000.00
Prepaid taxes: 1,000.00	Due to other banks and bankers: 1,000.00
Prepaid interest: 1,000.00	Notes and bills to be advanced: 1,000.00
Prepaid interest: 1,000.00	Other payables: 1,000.00
Prepaid interest: 1,000.00	Total: 2,000,000.00

(a) Excel+VBA

**Finding Error**

Row: 101, 101-01

Find Error

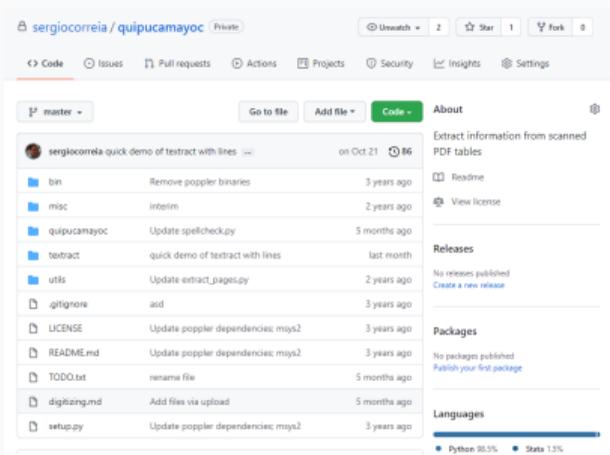
Error Type	Description	Status
1	URL	
2	Text	
3	Text	
4	Text	
5	Text	
6	Text	
7	Text	
8	Text	
9	Text	
10	Text	
11	Text	
12	Text	
13	Text	
14	Text	
15	Text	
16	Text	
17	Text	
18	Text	
19	Text	
20	Text	
21	Text	
22	Text	
23	Text	
24	Text	
25	Text	
26	Text	
27	Text	
28	Text	
29	Text	
30	Text	
31	Text	
32	Text	
33	Text	
34	Text	
35	Text	
36	Text	
37	Text	
38	Text	
39	Text	
40	Text	
41	Text	
42	Text	
43	Text	
44	Text	
45	Text	
46	Text	
47	Text	
48	Text	
49	Text	
50	Text	
51	Text	
52	Text	
53	Text	
54	Text	
55	Text	
56	Text	
57	Text	
58	Text	
59	Text	
60	Text	
61	Text	
62	Text	
63	Text	
64	Text	
65	Text	
66	Text	
67	Text	
68	Text	
69	Text	
70	Text	
71	Text	
72	Text	
73	Text	
74	Text	
75	Text	
76	Text	
77	Text	
78	Text	
79	Text	
80	Text	
81	Text	
82	Text	
83	Text	
84	Text	
85	Text	
86	Text	
87	Text	
88	Text	
89	Text	
90	Text	
91	Text	
92	Text	
93	Text	
94	Text	
95	Text	
96	Text	
97	Text	
98	Text	
99	Text	
100	Text	

Checklist not yet calculated

(b) Website

## Our tool - quipucamayoc

- Built in Python
- Modular tools implementing the methods we've used for these projects
- Should be easily combined with other packages (e.g. other layout parsers)



```
import quipucamayoc as q
```

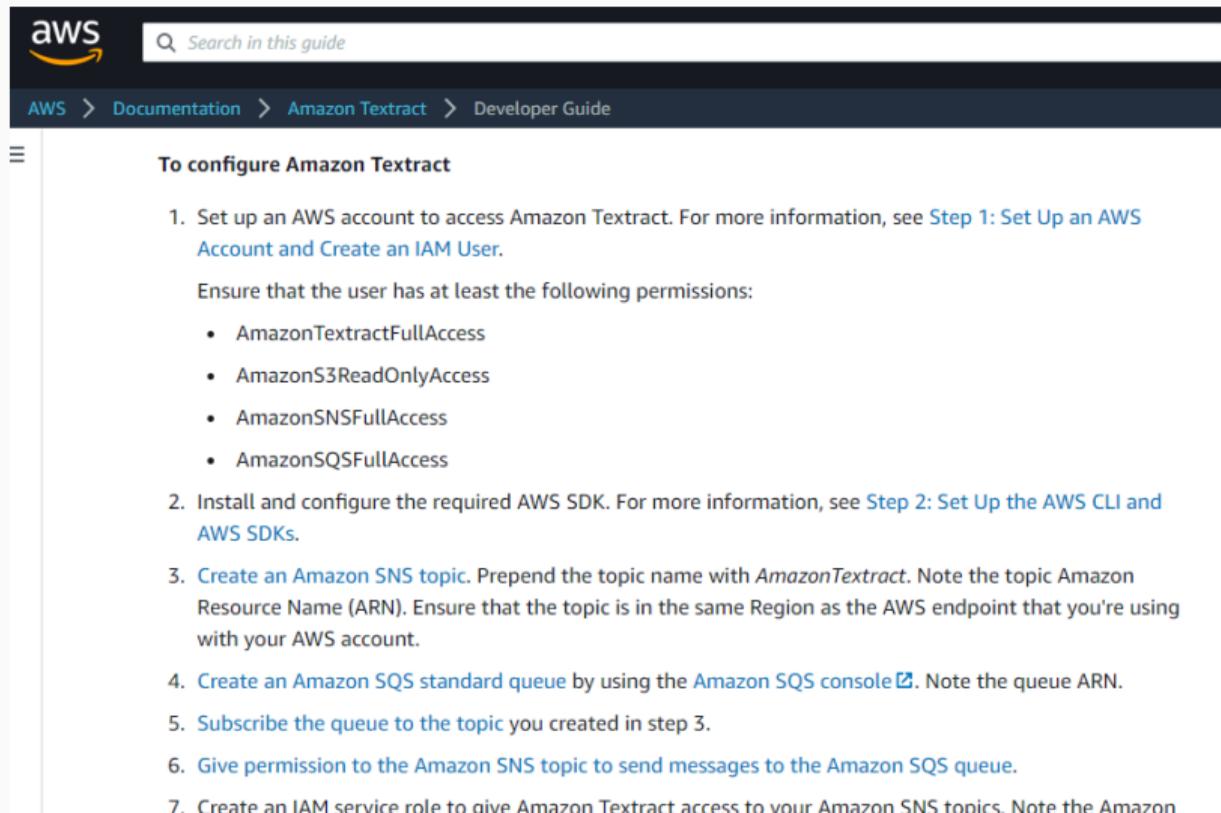
```
doc=q.read('rothschild-secret-ledger.pdf')  
page = doc.pages[0]
```

```
page.remove_black_background()  
page.deskew()  
page.binarize()
```

```
page.run_ocr(engine='amazon')  
page.run_ocr(engine='google')
```

```
s = q.spellchecker('frequencies.csv')  
process(page, spellchecker=s)
```

q.setup\_textextract() versus:



The screenshot shows the AWS documentation page for Amazon Text Extract. At the top left is the AWS logo. To its right is a search bar with the placeholder text "Search in this guide". Below the search bar is a breadcrumb trail: "AWS > Documentation > Amazon Text Extract > Developer Guide". On the left side of the page, there is a hamburger menu icon. The main content area is titled "To configure Amazon Text Extract" and contains a numbered list of steps:

1. Set up an AWS account to access Amazon Text Extract. For more information, see [Step 1: Set Up an AWS Account and Create an IAM User](#).  
Ensure that the user has at least the following permissions:
  - AmazonTextExtractFullAccess
  - AmazonS3ReadOnlyAccess
  - AmazonSNSFullAccess
  - AmazonSQSFullAccess
2. Install and configure the required AWS SDK. For more information, see [Step 2: Set Up the AWS CLI and AWS SDKs](#).
3. [Create an Amazon SNS topic](#). Prepend the topic name with *AmazonTextExtract*. Note the topic Amazon Resource Name (ARN). Ensure that the topic is in the same Region as the AWS endpoint that you're using with your AWS account.
4. [Create an Amazon SQS standard queue](#) by using the [Amazon SQS console](#). Note the queue ARN.
5. [Subscribe the queue to the topic you created in step 3](#).
6. [Give permission to the Amazon SNS topic to send messages to the Amazon SQS queue](#).
7. [Create an IAM service role to give Amazon Text Extract access to your Amazon SNS topics](#). Note the Amazon

1. Digitization at scale is a Leontief production function:
  - Only as good as its weakest step
  - Won't be successful if any of the steps (OCR, human review, data extraction, etc) is done poorly
  - Other researchers often mention poor results with OCR digitization; we suspect this is why

### 2. Use a cloud OCR provider:

- Not worth it to use a cheaper (free) OCR engine and then waste lots of time (=money) improving its results.
- Cloud OCR providers, today, are performant, robust, and still quite cheap (\$0.001-\$0.0015 per page)
- Cloud providers work in parallel without the need of us maintaining multiple servers.

### 3. Use the right tools when pre- and post-processing:

- Python worked better than {R, Stata}
- Well supported tools such as OpenCV worked better than state-of-the-art tools not yet battle-tested

### 4. Make human validation efficient:

- Human validation can be vastly sped up by using programs that help the reviewer:
  - Keyboard shortcuts
  - Side-to-side image-data comparisons
  - Auto zoom-in
  - Auto flag potential errors
- Not just about saving time, but about preserving *focus*

Thanks!