

## REPORT OF THE COMPTROLLER OF THE CURRENCY. 737

## NEW YORK.

## Chase National Bank, New York.

H. W. CANNON, *President.*

No. 2370.

J. T. MILLS, JR., *Cashier.*

Resources.		Liabilities.	
Loans and discounts.....	\$14,954,408.80	Capital stock paid in.....	\$500,000.00
Overdrafts.....	4,683.41	Surplus fund.....	1,000,000.00
U. S. bonds to secure circulation...	50,000.00	Undivided profits, less current expenses and taxes paid.....	284,769.70
U. S. bonds to secure deposits.....	200,000.00	National bank notes outstanding.....	45,000.00
U. S. bonds on hand.....	167,350.00	State bank notes outstanding.....	
Premiums on U. S. bonds.....	26,782.06	Due to other national banks.....	9,309,113.60
Stocks, securities, etc.....	2,882,266.11	Due to State banks and bankers..	9,312,523.30
Bank'g house, furniture, and fixtures		Dividends unpaid.....	
Other real estate and mortg'g's owned		Individual deposits.....	4,641,779.71
Due from other national banks.....	759,750.57	Certified checks.....	153,674.33
Due from State banks and bankers.	239,149.84	United States deposits.....	
Due from approved reserve agents.		Deposits of U.S. disbursing officers..	110,450.36
Checks and other cash items.....	5,855.93	Notes and bills rediscounted.....	
Exchanges for clearing house.....	375,878.15	Bills payable.....	
Bills of other national banks.....	45,250.00	Liabilities other than those above stated.....	
Fractional currency, nickels, cents.	262.13		
Specio.....	739,586.00		
Legal-tender notes.....	883,838.00		
U. S. certificates of deposit.....	4,020,000.00		
Redemption fund with Treas. U. S.	2,250.00		
Due from Treasurer U. S.....			
<b>Total.....</b>	<b>25,357,311.00</b>	<b>Total.....</b>	<b>25,357,311.00</b>

## OCR Technique: Naïve Output

REPORT OF DL A I 'THE COMPTROLLER OF THE CURRENCY.NE.  
Bank, Thomaston.Georges National Epwarp O'BRIEN  
U. S. certificates of deposit ..... mense e eae nce  
Bills PAY@DIGÂ® . 000020 see ce ek ces  
Due from U. S. Treasurer .....4,500 00  
\_-TOb@L 220000.245, 606 05000.00}  
245, 606 05

OCR on its own useless!

Three ways of improving OCR:

Multiple sources and engines

Topology of source

Autocorrection

# OCR Technique: Naïve Output

```
REPORT OF DL A I 'THE COMPTROLLER OF THE CURRENCY.NE.  
Bank, Thomaston.Georges National Epwarp O'BRIEN  
U. S. certificates of deposit ..... mense e eae nce  
Bills PAY@DIGÂ® . 000020 see ce ek ces  
Due from U. S. Treasurer .....4,500 00  
_-TOB@L 220000.245, 606 05000.00}  
245, 606 05
```

→ OCR on its own useless!

Three ways of improving OCR:

1. Multiple sources and engines
2. Topology of source
3. Autocorrection

## OCR Improvement 1: Ensemble of Scans and Engines

- Old documents often have stains, ink blots and broken pages
- Different OCR engines make different errors: one confuses *BOB* with *808*; another confuses *bonds* with *bomls*

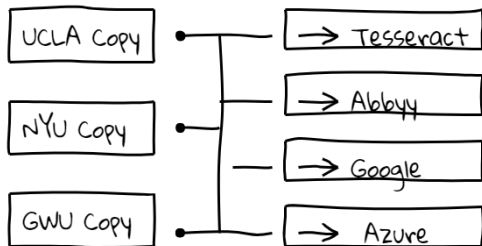
### **Second National Bank, Washington.**

<i>it.</i>	No. 2038.	
.....	\$252, 139 12	Capital stock paid in.....
.....	4, 503 20	
tion .....	170, 000 00	Surplus fund.....
:s.....	.....	Undivided profits.....
rtgages ..	55, 820 8	National bank notes outst:
.....	9, 986 46	State bank notes outstand

[back](#)

## OCR Improvement 1: Ensemble of Scans and Engines

- We gather as many different documents as possible and apply several modern OCR engines
- For example, in 1895 there are 12 inputs from three copies and four OCR engines:



[back](#)

# OCR Improvement 2: Use Computer Vision to Detect Columns and Rows

- To separate each column, we use OpenCV to detect all lines and separate columns (Konrad 2017)

National Exchange Bank, Augusta.			
ALFRED BAKER, <i>President.</i>		JOHN CRAIG, <i>Cashier.</i>	
No. 4860.			
Resources.		Liabilities.	
Loans and discounts .....	\$275,419 19	Capital stock .....	\$300,000 00
Overdrafts .....		Surplus fund .....	3,000 00
U. S. bonds to secure circulation .....	300,000 00	Undivided profits .....	14,548 89
U. S. bonds to secure deposits .....		National bank notes outstanding .....	270,000 00
U. S. bonds and securities on hand .....		State bank notes outstanding .....	
Other stocks, bonds, and mortgages .....	2,550 00	Dividends unpaid .....	2,610 00
Due from redeeming agents .....	35,412 23	Individual deposits .....	147,798 54
Due from other national banks .....	129 37	U. S. deposits .....	
Due from State banks and bankers .....	16,383 95	Deposits of U. S. disbursing officers .....	
Real estate, furniture, and fixtures .....	2,483 00	Due to national banks .....	1,000 00
Current expenses .....	3,123 54	Due to State banks and bankers .....	
Premiums paid .....	31,500 00	Notes and bills re-discounted .....	
Checks and other cash items .....	24,866 84	Bills payable .....	
Exchanges for clearing house .....			
Bills of other national banks .....	9,731 00		
Fractional currency .....	55 31		
Specie .....			
Legal tender notes .....	37,300 00		
Three per cent. certificates .....			
<b>Total.....</b>	<b>738,957 43</b>	<b>Total.....</b>	<b>738,957 43</b>

back

## OCR Improvement 3: Autocorrect Errors

- Construct list of possible balance sheet items, city names, etc.
- Feed to a spelling corrector (Norvig 2016)
- **Example:** Washlngion → Washington
- For numbers, enforce heuristics (no leading 0, three digits between commas, etc.)
- Allow OCR engines to vote, digit by digit:

True value:	123456	
Abbyy OCR:	.23456	Skipped first digit
Tesseract OCR:	128450	Confused 3 and 8
Azure OCR:	123156	Confused 1 and 4
Select common digit:	123456	

## OCR Improvement 3: Autocorrect Errors

- Construct list of possible balance sheet items, city names, etc.
- Feed to a spelling corrector (Norvig 2016)
- **Example:** Washlngion → Washington
- For numbers, enforce heuristics (no leading 0, three digits between commas, etc.)
- Allow OCR engines to vote, digit by digit:

---

True value:	123456	
Abbyy OCR:	_23456	Skipped first digit
Tesseract OCR:	128450	Confused 3 and 8
Azure OCR:	123156	Confused 1 and 4
Select common digit:	123456	

---



# Data

- Accurately identified all banks, including charter numbers
  - Geolocated cities (Schmidt 2017)
  - Major events: entry, exit (receivership, liquidation, mergers)
- Assess accuracy through balance sheet identities:
  - Total Assets = Total Liabilities (97% of cases)
  - Total Assets =  $\sum \text{Assets}_i$  (65%)
  - Total Liabilities =  $\sum \text{Liabilities}_i$  (80%)
  - **All three: 52% of cases.**
- Residual gets fixed by hand.

back