



# Automated Digitization of the Censuses of Housing Block Statistics, 1940-1970

Jeffrey Lin, Dan Moulton, Isaac Rand & Robyn Smith  
Federal Reserve Bank of Philadelphia

---

August 2024

PhiladelphiaFed.org | @PhiladelphiaFed



## Disclaimer

The views expressed here are those of the authors and do not necessarily represent the views of the Federal Reserve Bank of Philadelphia or the Federal Reserve System.

# Digitizing Block Statistics

- What
- Why
- Goals
- Tasks and Challenges



# Digitizing Block Statistics

- **What**
- Why
- Goals
- Tasks and Challenges

Census of Housing  
Block Statistics 

HOUSING—BLOCK STATISTICS

Table 3.—CHARACTERISTICS OF HOUSING FOR CENSUS TRACTS BY BLOCKS: 1940—Con.

Census tract	Block	Total structures	ALL DWELLING UNITS BY OCCUPANCY AND TENURE					ALL DWELLING UNITS BY YEAR BUILT				OCCUPIED DWELLING UNITS			ALL DWELLING UNITS BY STATE OF REPAIR AND PLUMBING EQUIPMENT				OWNER-OCCUPIED UNITS BY MORTGAGE STATUS		ALL DWELLING UNITS BY CONTRACT OR ESTIMATED RENT	
			Total dwelling units	Owner occupied	Tenant occupied	Vacant, for sale or rent	Vacant, other	Number reporting	1930 to 1940	1920 to 1929	1900 to 1919	1899 or before	Total occupied	Occupied by non-white	Persons per room Num- ber or rptg. more	Number reporting	Needing repair or no private bath	Need- ing re- pair	No pri- vate bath	Number reporting	Mort- gaged	Number reporting
3-A	24	21	34	11	17	6	34				34	28	28	2	19			6	2	33	24.24	
	25	41	43	20	20	3	42				42	40	39	1	20			18	9	43	18.21	
	26	32	36	16	19	1	36	2	4		36	35	35	5	17			14	11	36	18.22	
	27	34	38	13	24	1	37		4		37	37	37	1	15	1	15	11	6	37	23.73	
	28	26	49	18	30	1	49	1	3		48	48	48	7	24			24	18	42	21.92	
	29	18	24	12	12		24	1	2		24	24	24	2	8		8	12	6	22	22.55	
	30	11	25	4	20	1	25				24	24	24	2	15		15	4	2	22	17.86	
	31	24	38	12	23	3	38				35	35	35	3	35	35	15	12	5	38	20.66	
	32	18	33	11	20	2	32				31	31	28	3	23		23	10	6	31	20.42	
	33	28	34	6	27		34		3		29	33	33	5	29		29	5	5	34	18.71	
	34	6	9	3	6		9				9	9	9	2	6		6	3	1	9	20.00	
	35	28	30	13	16	1	30				28	29	29	1	19		3	19	3	30	21.57	
	36	22	31	7	23	1	31				31	30	30	4	18		10	3	2	31	17.32	
	37	23	43	7	32	4	43	4			39	39	39	3	24		15	20	4	42	22.14	
	38	20	26	4	20	2	26				24	24	24	4	17		10	11	1	26	18.04	
	39	41	44	12	28	4	44	1			40	40	40	1	14		14	4	4	44	28.32	
	40	43	71	21	47	3	71				68	68	68	4	25		1	24	15	71	24.70	
	41	27	36	10	25	1	36				35	35	35	6	27		25	20	10	36	17.33	
	42	28	50	4	43	3	49				49	47	46	1	29		22	4	2	49	17.80	
	43	2	3	3	3		3				3	3	3	1	3		3			3	18.67	

# Census of Housing Block Statistics

- Most granular, earliest, extant Census spatial data on housing.
- 1940-1970.
- Tens of thousands+ of scanned pages of tables and maps.



## What's in it?

- Tenure, occupancy, structure age and condition, rents and values, race of occupants.
- All houses, not just occupied ones.
- High level of spatial detail: Usually, a city block.
- Small size (Pop. ~50 vs ~4,000 for ED/Tract).
- Coverage of large section of cities over time.
- 191 cities in 1940 → All 1970 urbanized areas.

## What's it good for?

Studies of housing investment and maintenance and long-run urban dynamics.

Studies of policies and processes that occur at extremely localized spatial scales.

Studies of many cities, or a single city's history.

# Digitizing Block Statistics

- What
- **Why**
- Goals
- Tasks and Challenges

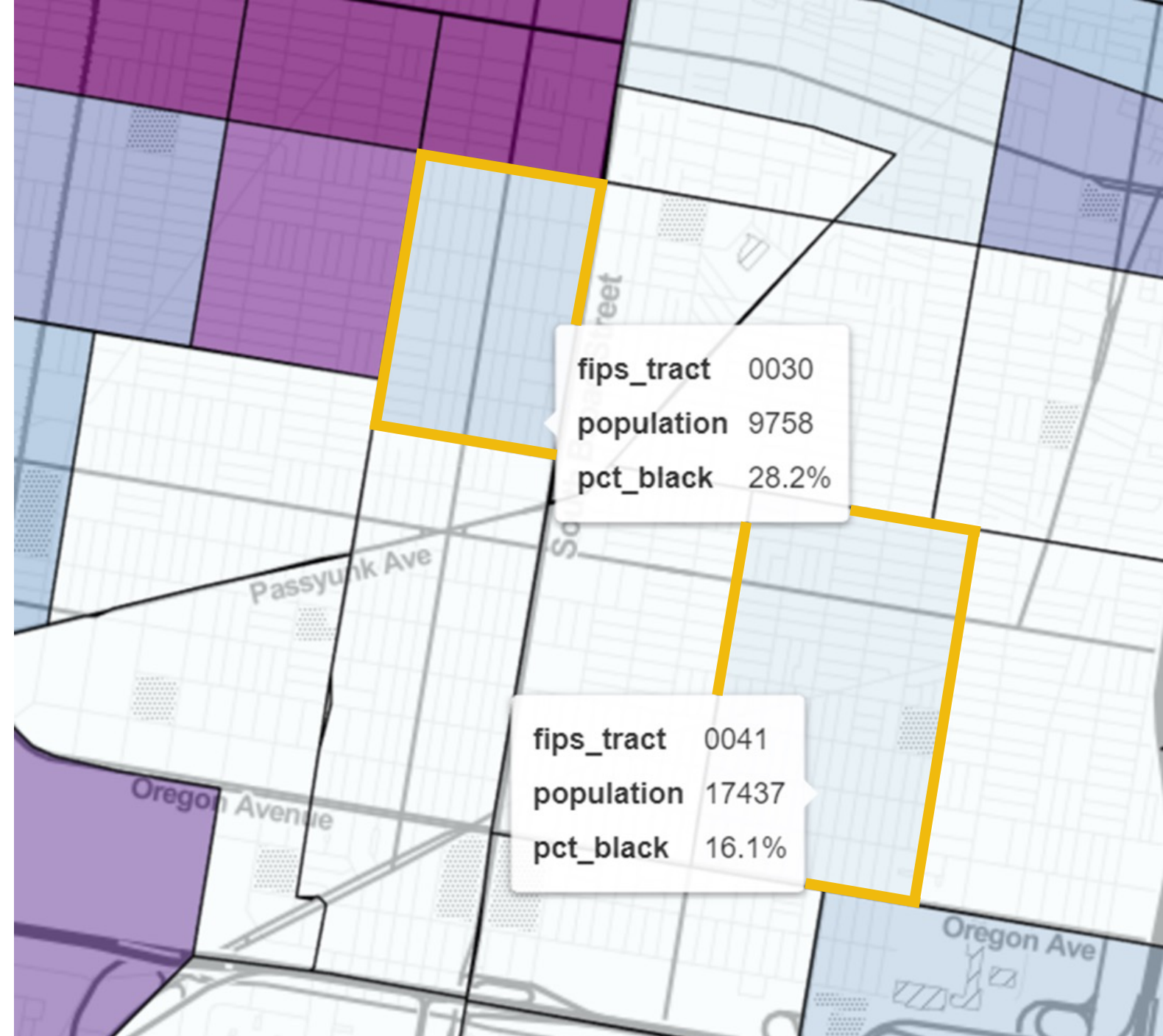
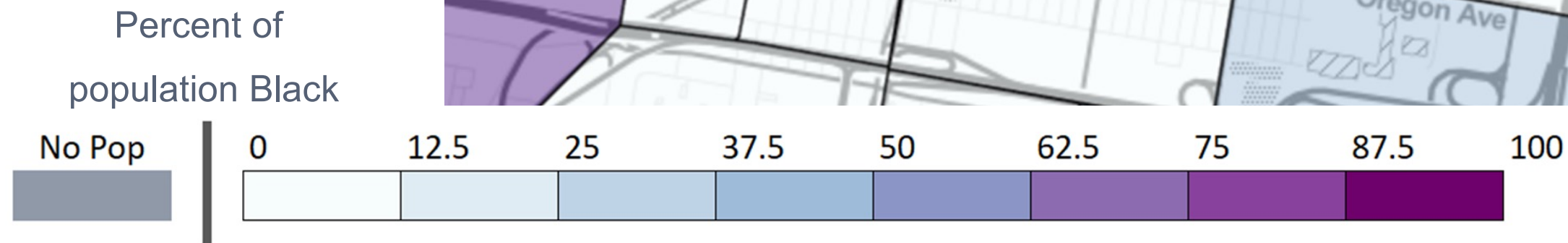
## EXAMPLES

Localized processes

Localized policies

# Localized Processes

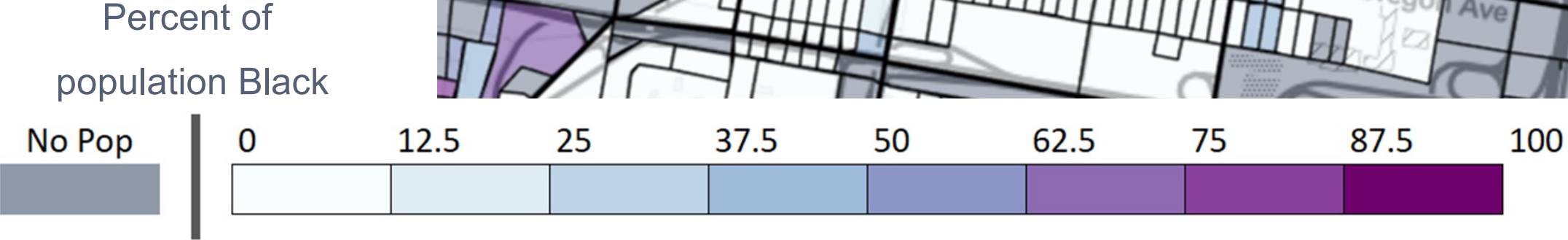
- Residential segregation in South Philadelphia, 1970 at **Census Tract** scale





# Localized Processes

- Residential segregation in South Philadelphia, 1970 at **Census Block** scale

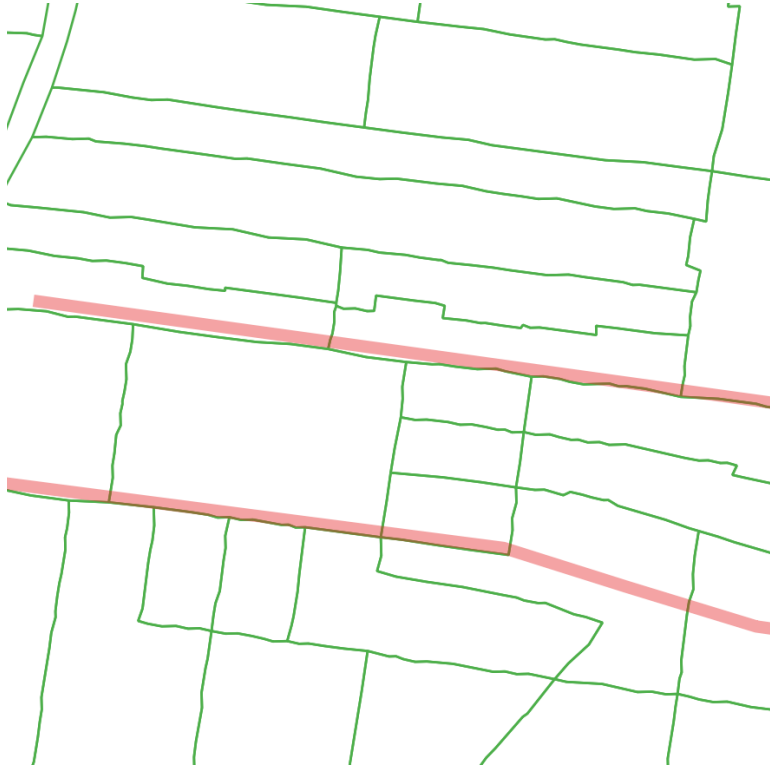


# Localized Policies

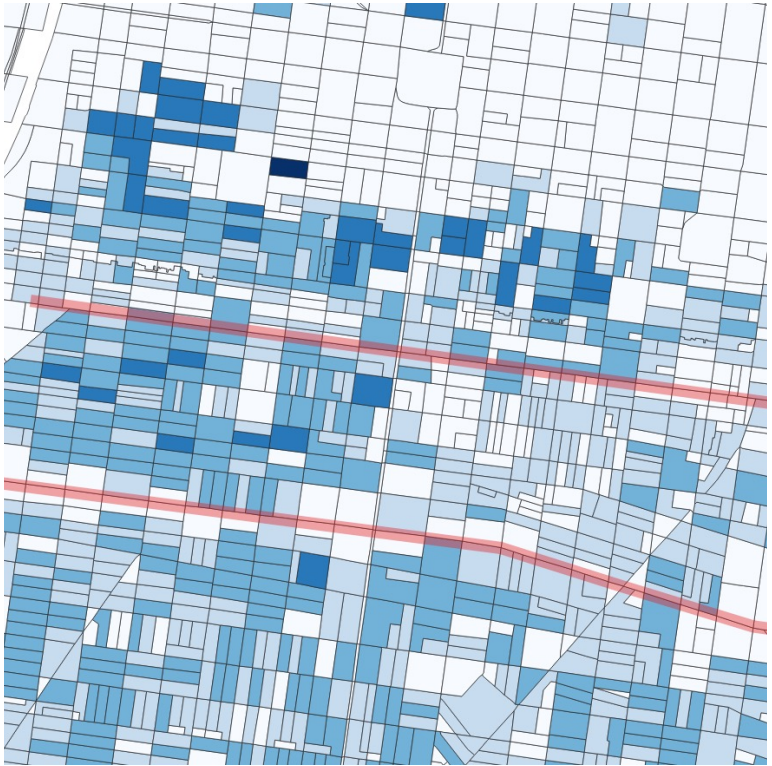
## Runner-up design

- “Expecting an Expressway” (Brinkman, Lin & Mangum).
- Two proposed routes for the Crosstown Expressway in South Philadelphia.

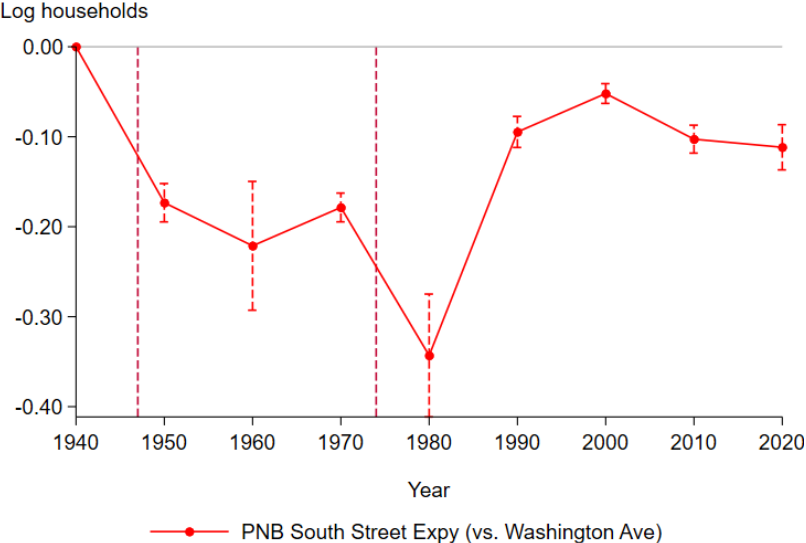
Tracts



Blocks



Difference in differences



# Digitizing Block Statistics

- What
- Why
- **Goals**
- Tasks and Challenges



# Our Goals



- Block data for **16 cities**, 1940-1970.
- Training and validation data.
- Code and methods.
- **Freely distributed for use and re-use.**

# Digitizing Block Statistics

- What
- Why
- Goals
- **Tasks and Challenges**

## Three Tasks

Shapes  
Situations  
Statistics

## Challenges

Limitations of  
traditional  
approaches  
Our current work

# 3 Tasks, 3 Pieces of Data

1

Shapes

Blocks need to know their:

2

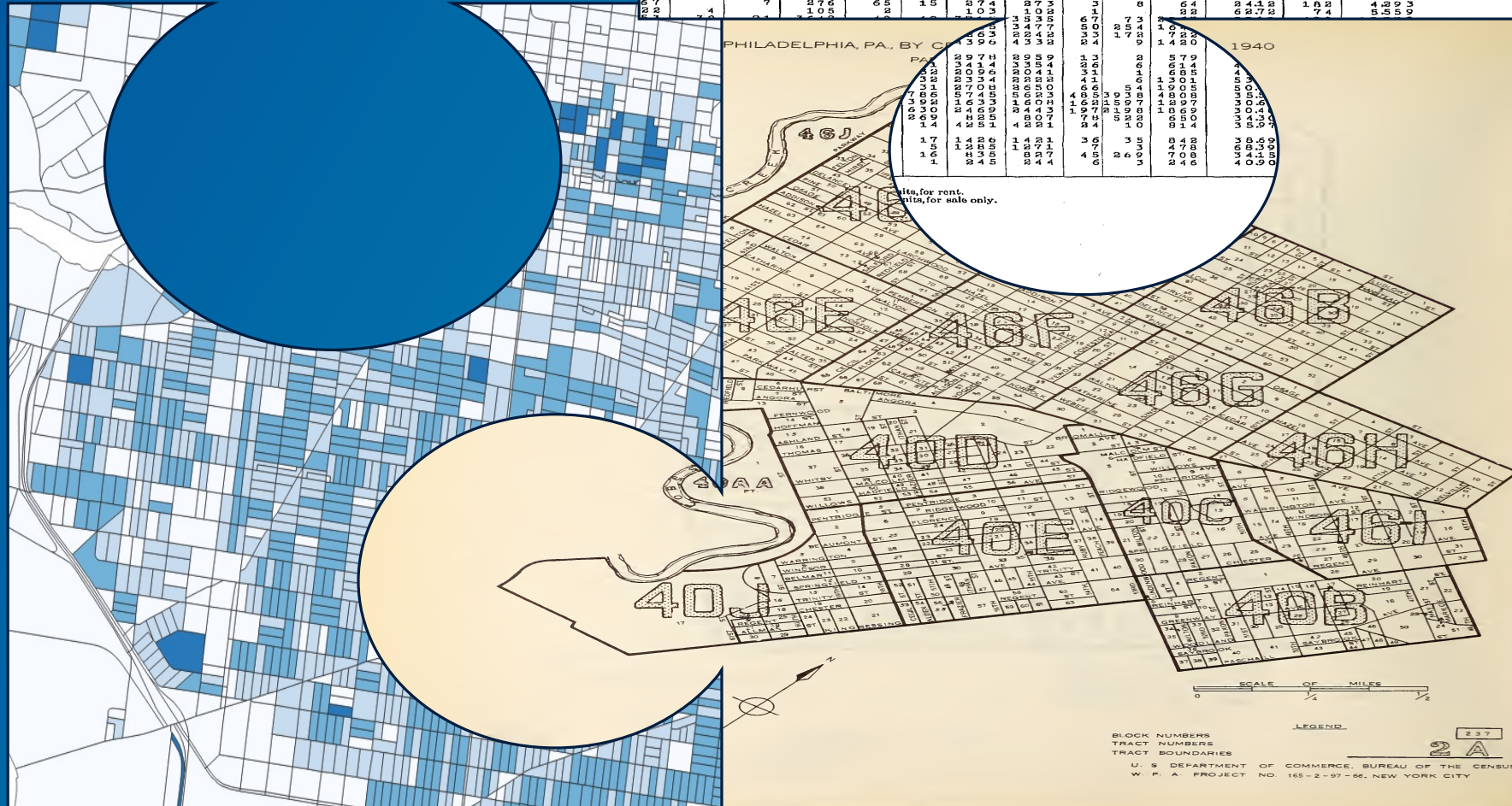
Situations

3

Statistics

2.—CHARACTERISTICS OF HOUSING BY CENSUS TRACTS: 1950—Con.

Tract	Area (sq. ft.)	Population	All dwelling units by condition and plumbing facilities			Occupied dwelling units			Contract monthly rent <sup>1</sup>		Value <sup>2</sup> of one-dwelling-unit structures		
			Number reporting	No private bath or dilap.	No running water or dilap.	Total	Persons per room	1.01 or more	Occupied by non-whites	Number reporting	Average monthly rent (dollars)	Number reporting	Average value (dollars)
4000	...	...	...	...	...	...	...	...	...	...	...	...	
4001	...	...	...	...	...	...	...	...	...	...	...	...	
4002	...	...	...	...	...	...	...	...	...	...	...	...	
4003	...	...	...	...	...	...	...	...	...	...	...	...	
4004	...	...	...	...	...	...	...	...	...	...	...	...	
4005	...	...	...	...	...	...	...	...	...	...	...	...	
4006	...	...	...	...	...	...	...	...	...	...	...	...	
4007	...	...	...	...	...	...	...	...	...	...	...	...	
4008	...	...	...	...	...	...	...	...	...	...	...	...	
4009	...	...	...	...	...	...	...	...	...	...	...	...	
4010	...	...	...	...	...	...	...	...	...	...	...	...	
4011	...	...	...	...	...	...	...	...	...	...	...	...	
4012	...	...	...	...	...	...	...	...	...	...	...	...	
4013	...	...	...	...	...	...	...	...	...	...	...	...	
4014	...	...	...	...	...	...	...	...	...	...	...	...	
4015	...	...	...	...	...	...	...	...	...	...	...	...	
4016	...	...	...	...	...	...	...	...	...	...	...	...	
4017	...	...	...	...	...	...	...	...	...	...	...	...	
4018	...	...	...	...	...	...	...	...	...	...	...	...	
4019	...	...	...	...	...	...	...	...	...	...	...	...	
4020	...	...	...	...	...	...	...	...	...	...	...	...	
4021	...	...	...	...	...	...	...	...	...	...	...	...	
4022	...	...	...	...	...	...	...	...	...	...	...	...	
4023	...	...	...	...	...	...	...	...	...	...	...	...	
4024	...	...	...	...	...	...	...	...	...	...	...	...	
4025	...	...	...	...	...	...	...	...	...	...	...	...	
4026	...	...	...	...	...	...	...	...	...	...	...	...	
4027	...	...	...	...	...	...	...	...	...	...	...	...	
4028	...	...	...	...	...	...	...	...	...	...	...	...	
4029	...	...	...	...	...	...	...	...	...	...	...	...	
4030	...	...	...	...	...	...	...	...	...	...	...	...	
4031	...	...	...	...	...	...	...	...	...	...	...	...	
4032	...	...	...	...	...	...	...	...	...	...	...	...	
4033	...	...	...	...	...	...	...	...	...	...	...	...	
4034	...	...	...	...	...	...	...	...	...	...	...	...	
4035	...	...	...	...	...	...	...	...	...	...	...	...	
4036	...	...	...	...	...	...	...	...	...	...	...	...	
4037	...	...	...	...	...	...	...	...	...	...	...	...	
4038	...	...	...	...	...	...	...	...	...	...	...	...	
4039	...	...	...	...	...	...	...	...	...	...	...	...	
4040	...	...	...	...	...	...	...	...	...	...	...	...	



# 3 Tasks, 3 Pieces of Data

1

**Shape**

Segmenting Block Shapes from Maps

2

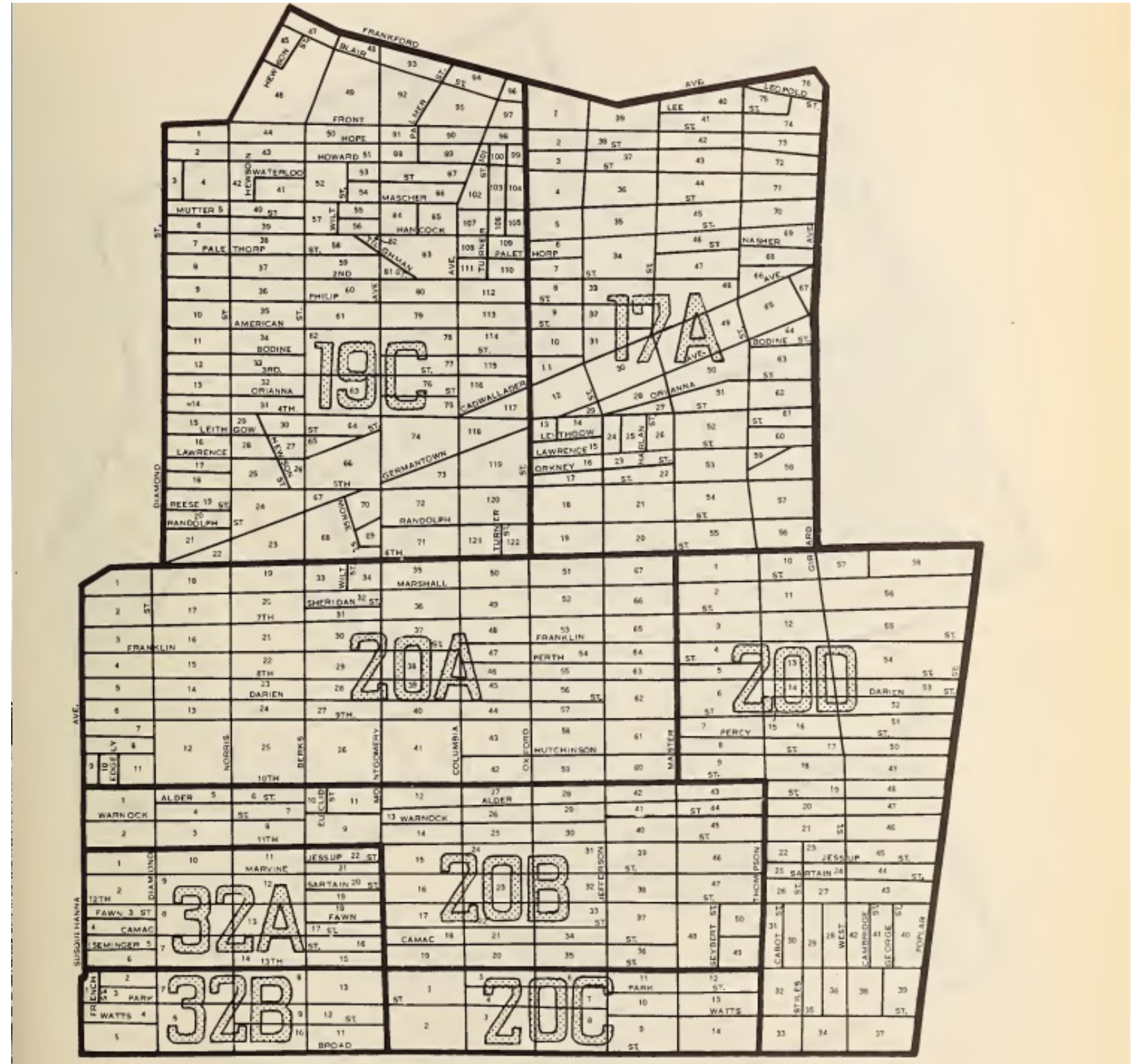
Situation

3

Statistics

# Our Ideal Process Has Only Three Steps

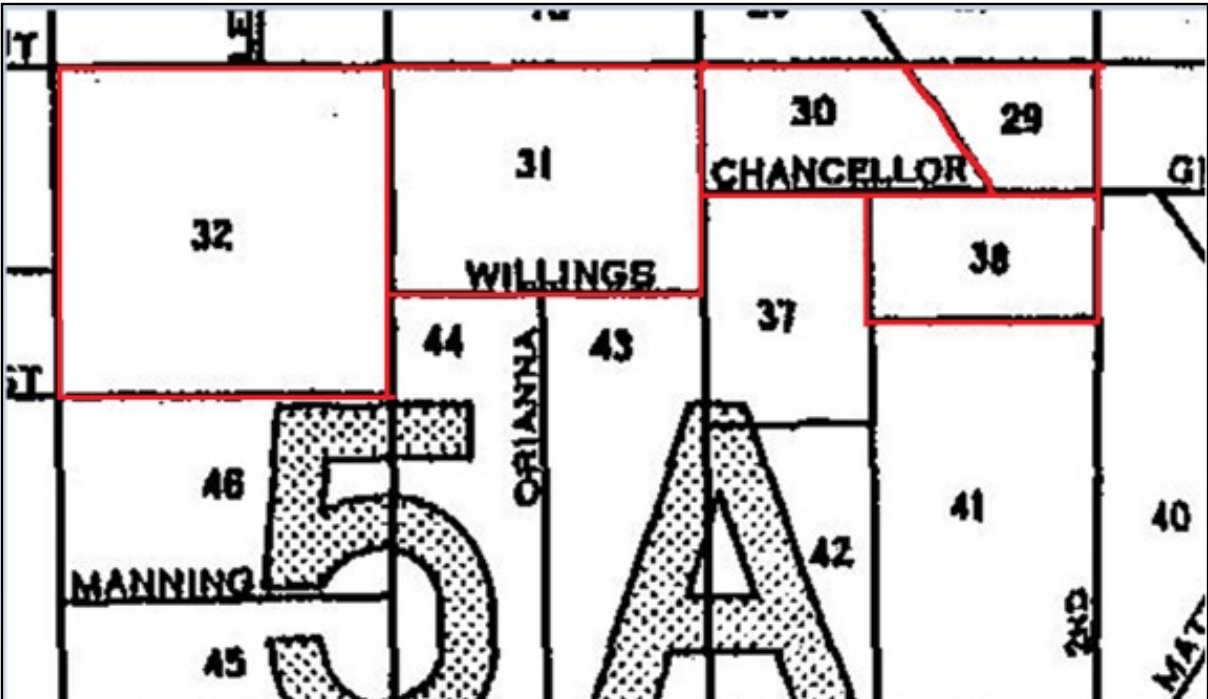
1. Identify closed loops of black ink.
2. Call them all blocks.
3. Declare victory.





# Unfortunately, This Process Fails Spectacularly

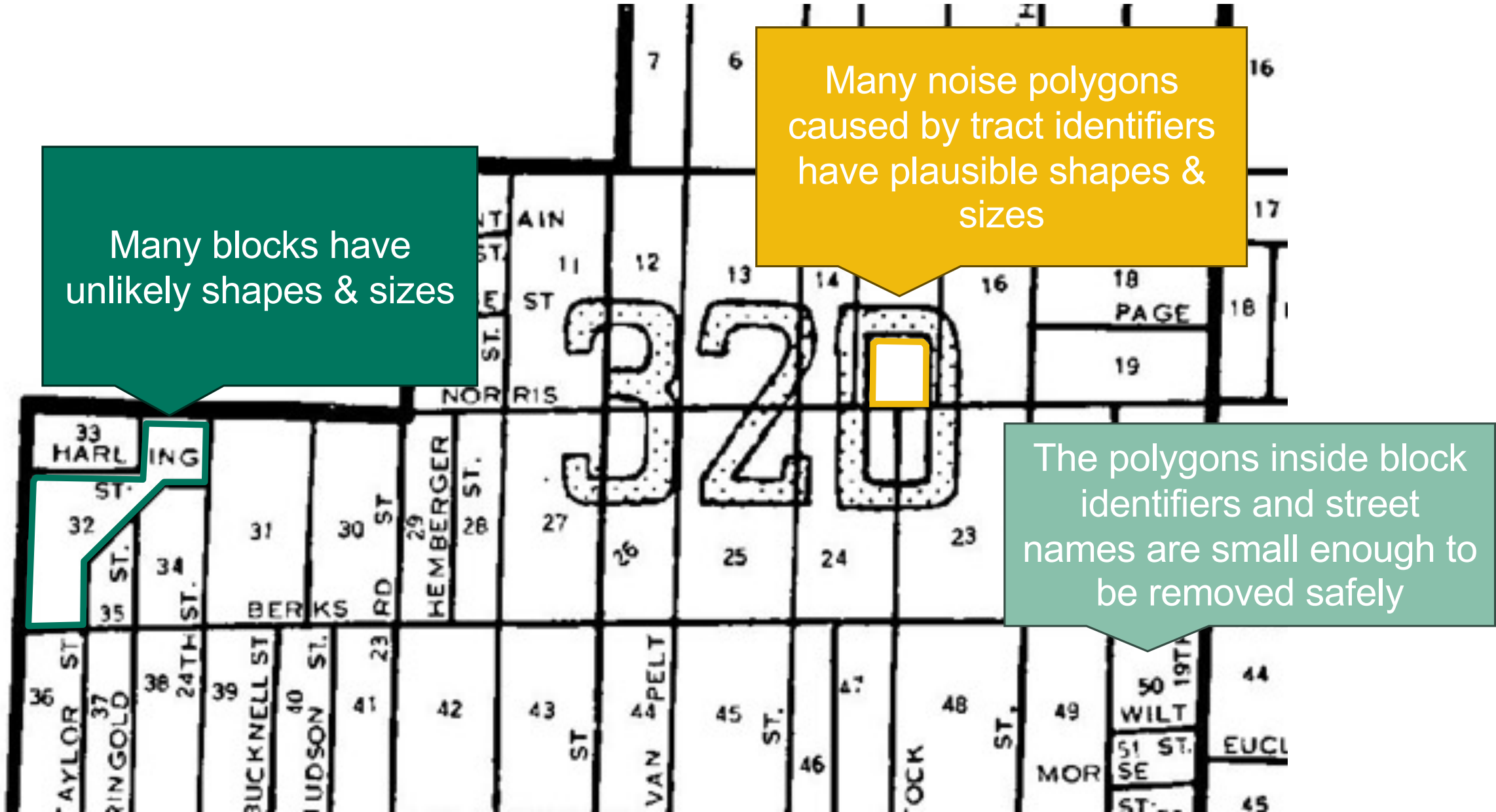
While many closed loops of black ink are blocks....



Many closed loops of black ink are not blocks 😞

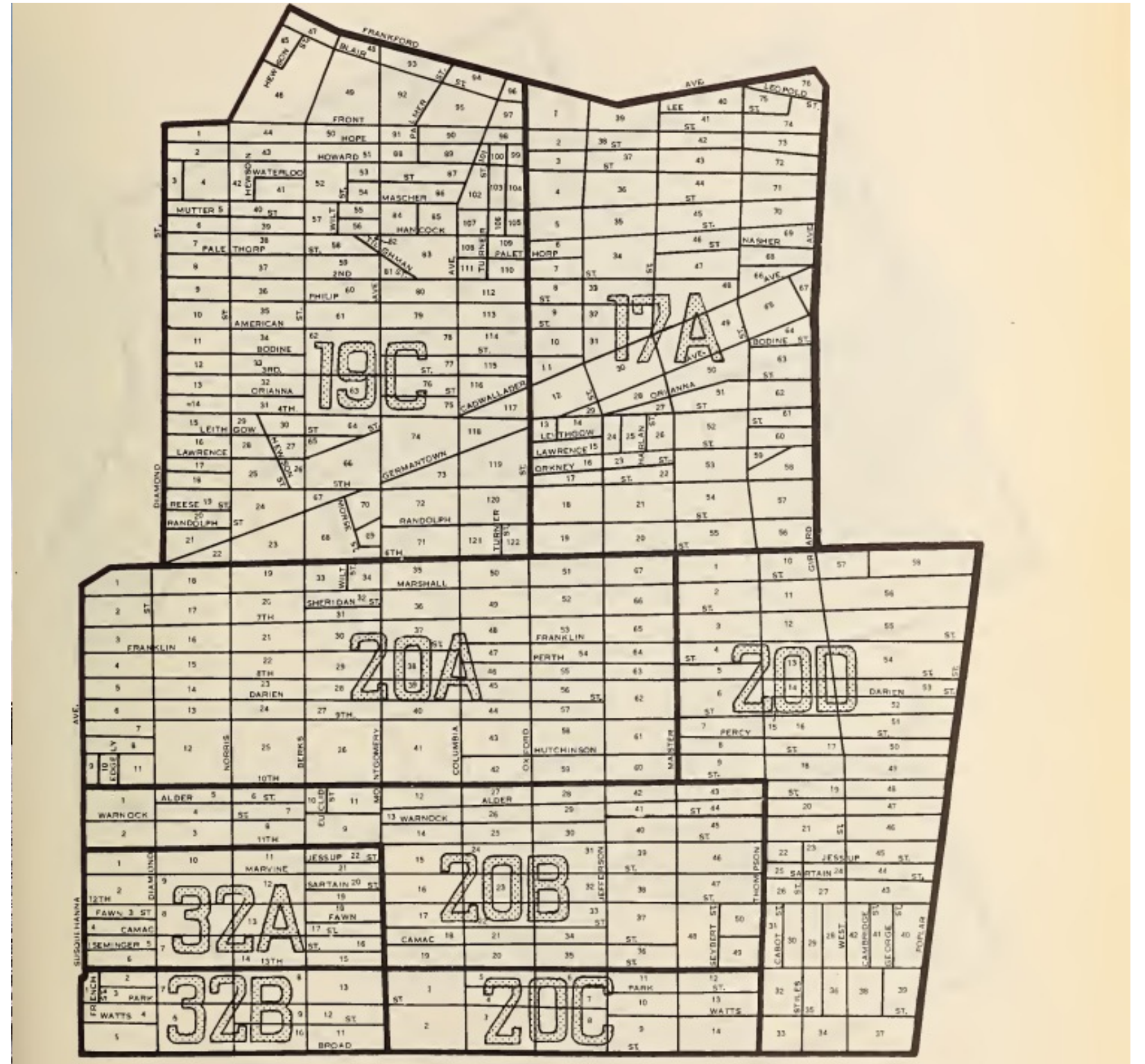


# Can We Handle the Noise?



# Our Ideal Process Has Only ~~Three~~ Four Steps

1. Remove the tract identifiers from the page
2. Identify remaining closed loops of black ink
3. Call any reasonably large loops blocks
4. Declare victory



How can we remove  
tract identifiers?

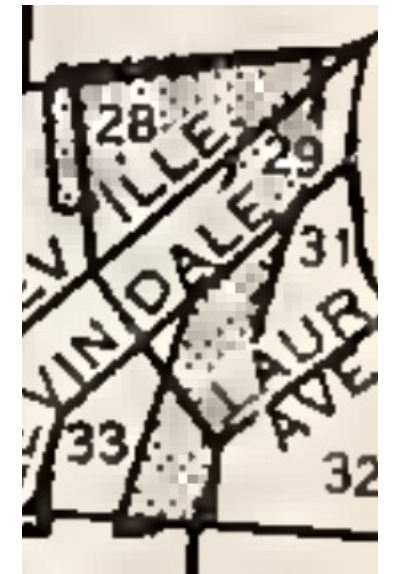
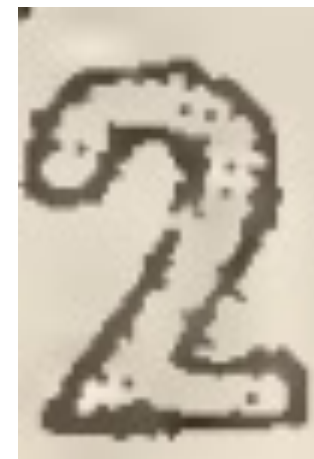
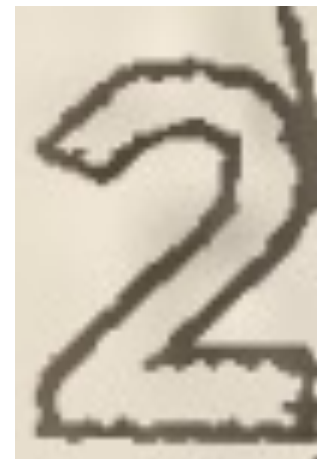
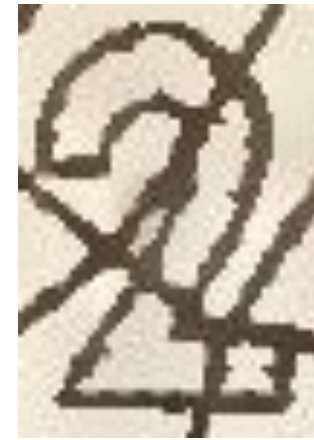
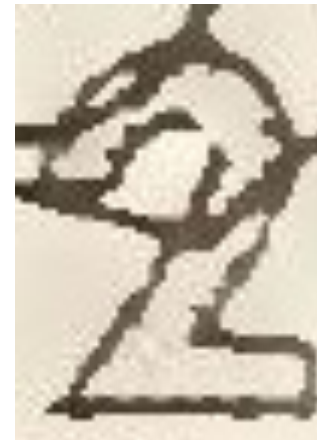
# Traditional Method 1: Matching Large Shapes

## Issues

- Arbitrary Rotations, Inconsistent Scale, Shape, and Font, Noise/Interference from other features.

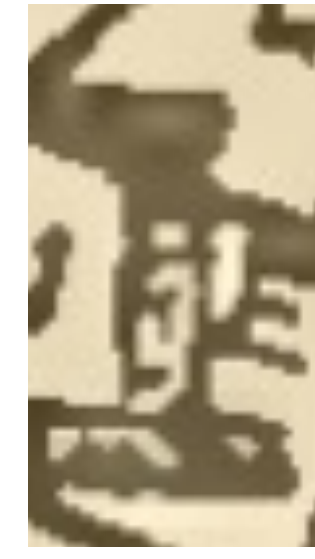
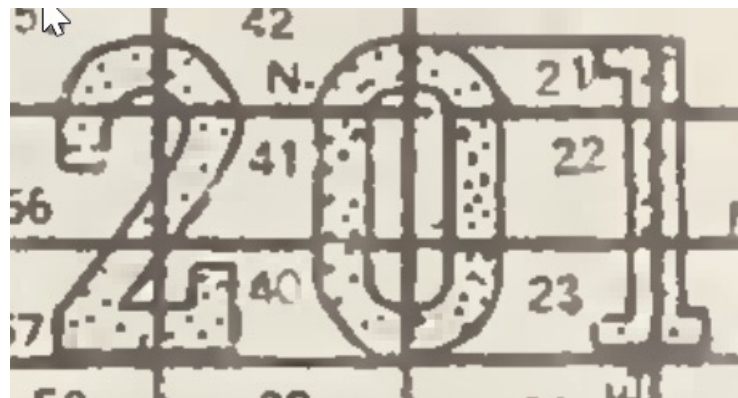
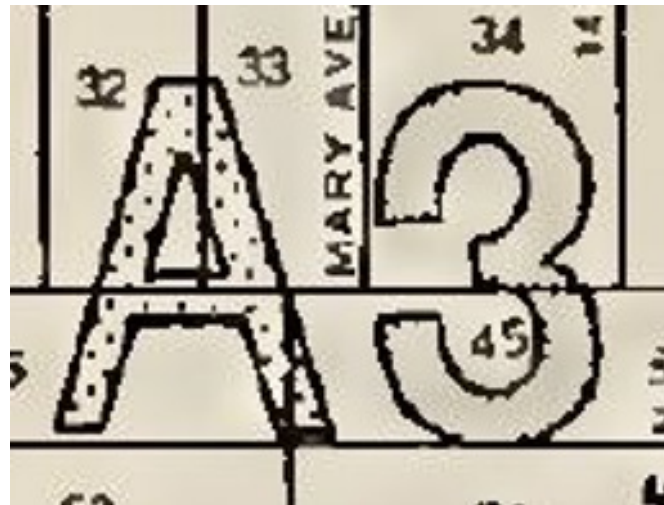
## Low-confidence matches

- If we accept low confidence matches of enough templates, everything starts to look like a tract identifier.
- Especially problematic with blocky characters like 1 and E.



# Traditional Method 2: Matching Patterns

- Sometimes speckled
- Sometimes no fill
- Block boundaries still a problem
- Sometimes inked



# Current Work: CNN

## More holistic

- Consider properties of block boundaries as well as properties of tract identifiers.
- Focus on identifying block boundaries, not removing tract identifiers.

## More flexible

- Can learn more patterns than we can with shape template matching.
- Can address partial shapes.
- Important because of intersections between boundaries and identifiers.

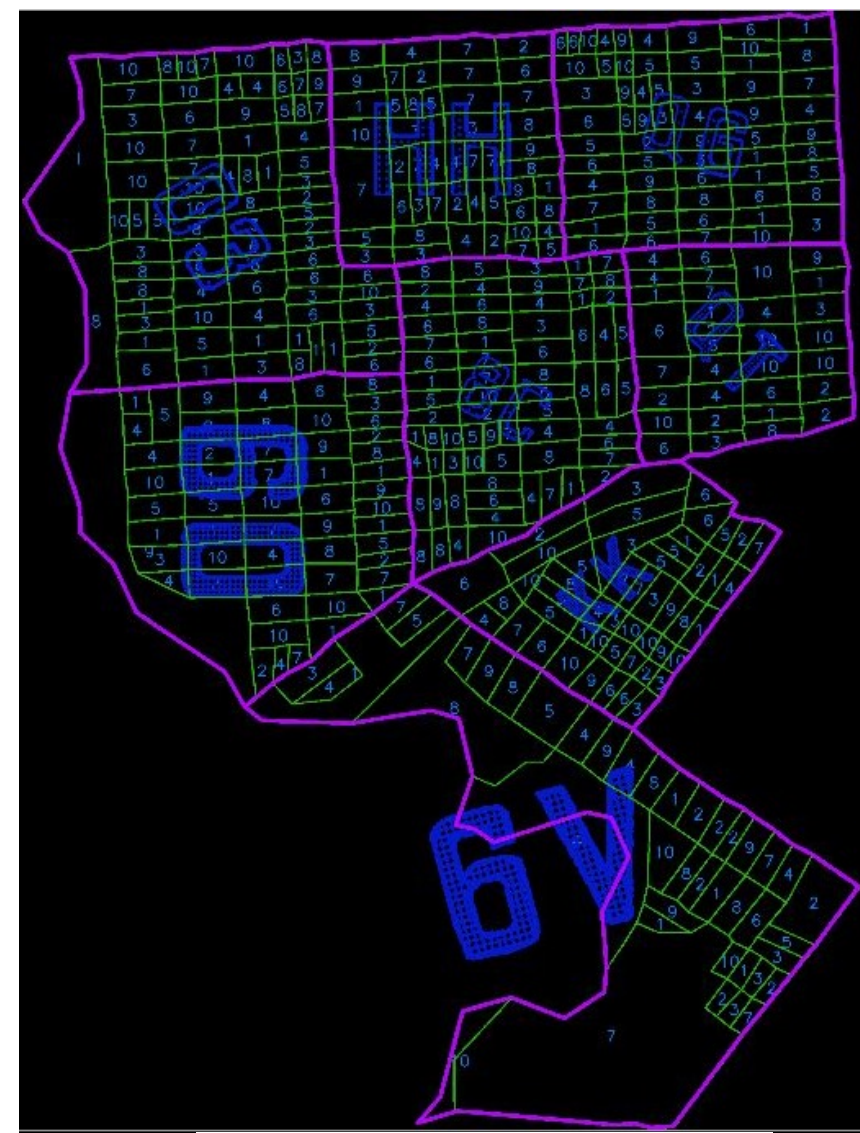


# Creating Training Data

- Hand annotations are expensive; Simulating maps is cheap.
- Sample 1990 Census block and tract boundaries from NHGIS.
- Sample tract and block identifiers from real 1940 maps.
- Randomly assign speckle density to tract identifiers.



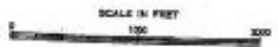
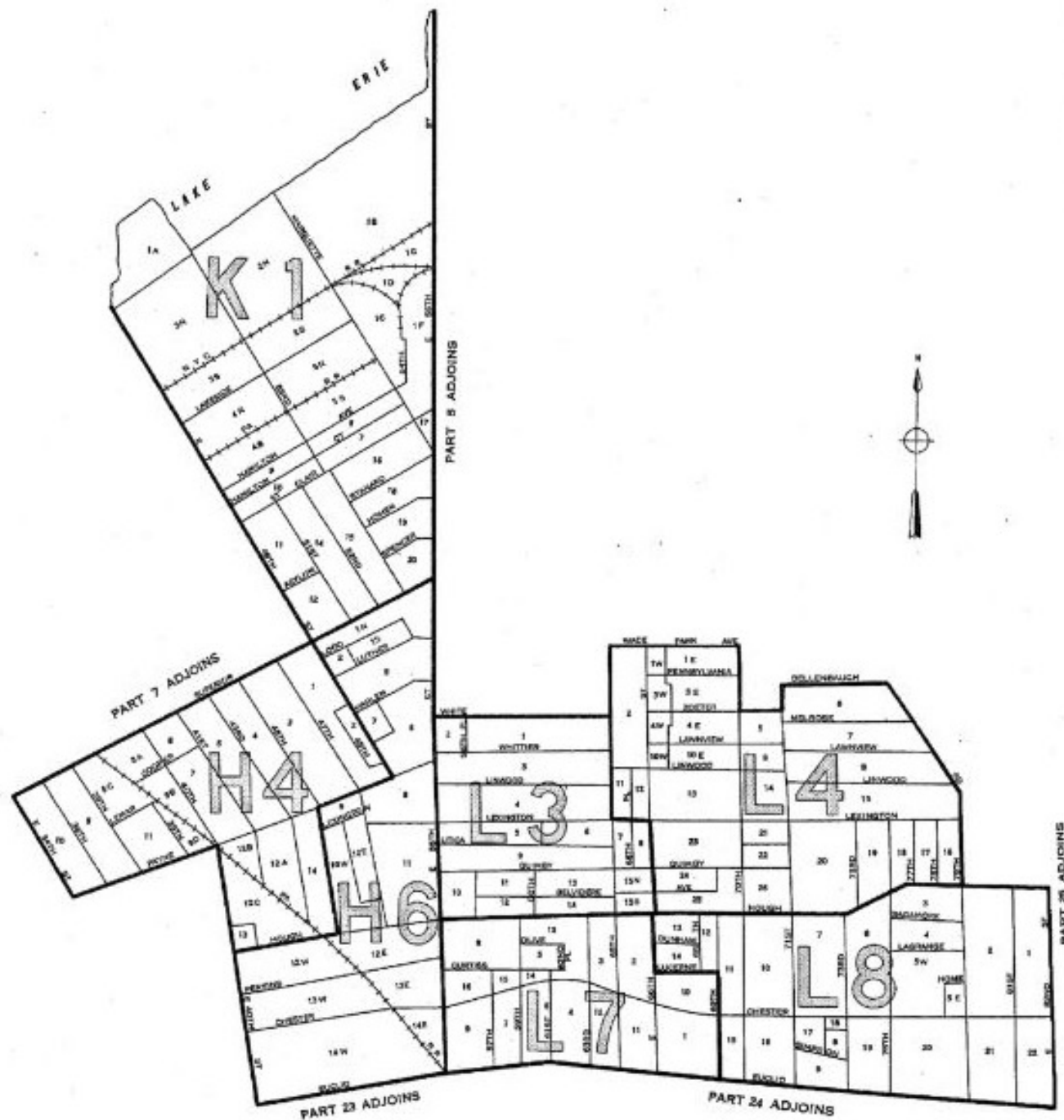
Simulated Map



Training Mask



Can the model trained on simulated maps generalize to real ones?



LEGEND

BLOCK NUMBERS  
TRACT NUMBERS  
TRACT BOUNDARIES

U.S. DEPARTMENT OF COMMERCE, BUREAU OF THE CENSUS

25  
2



LEGEND

BLOCK NUMBERS  
TRACT NUMBERS  
TRACT BOUNDARIES

U.S. DEPARTMENT OF COMMERCE, BUREAU OF THE CENSUS

25

2



SCALE OF MILES  
 0 1/4 1/2

LEGEND  
 BLOCK NUMBERS  
 TRACT NUMBERS  
 TRACT BOUNDARIES



SCALE OF MILES  
 0 1/4 1/2

LEGEND  
 BLOCK NUMBERS  
 TRACT NUMBERS  
 TRACT BOUNDARIES

# What's Next?

## Model and training improvements

- Better simulated maps.
- Augment with hand annotations.

## Add more steps

- Inpainting lines erased by CNN.
- Suggestions?



# And Now For Something Completely Different

(1970 maps)

## Promises

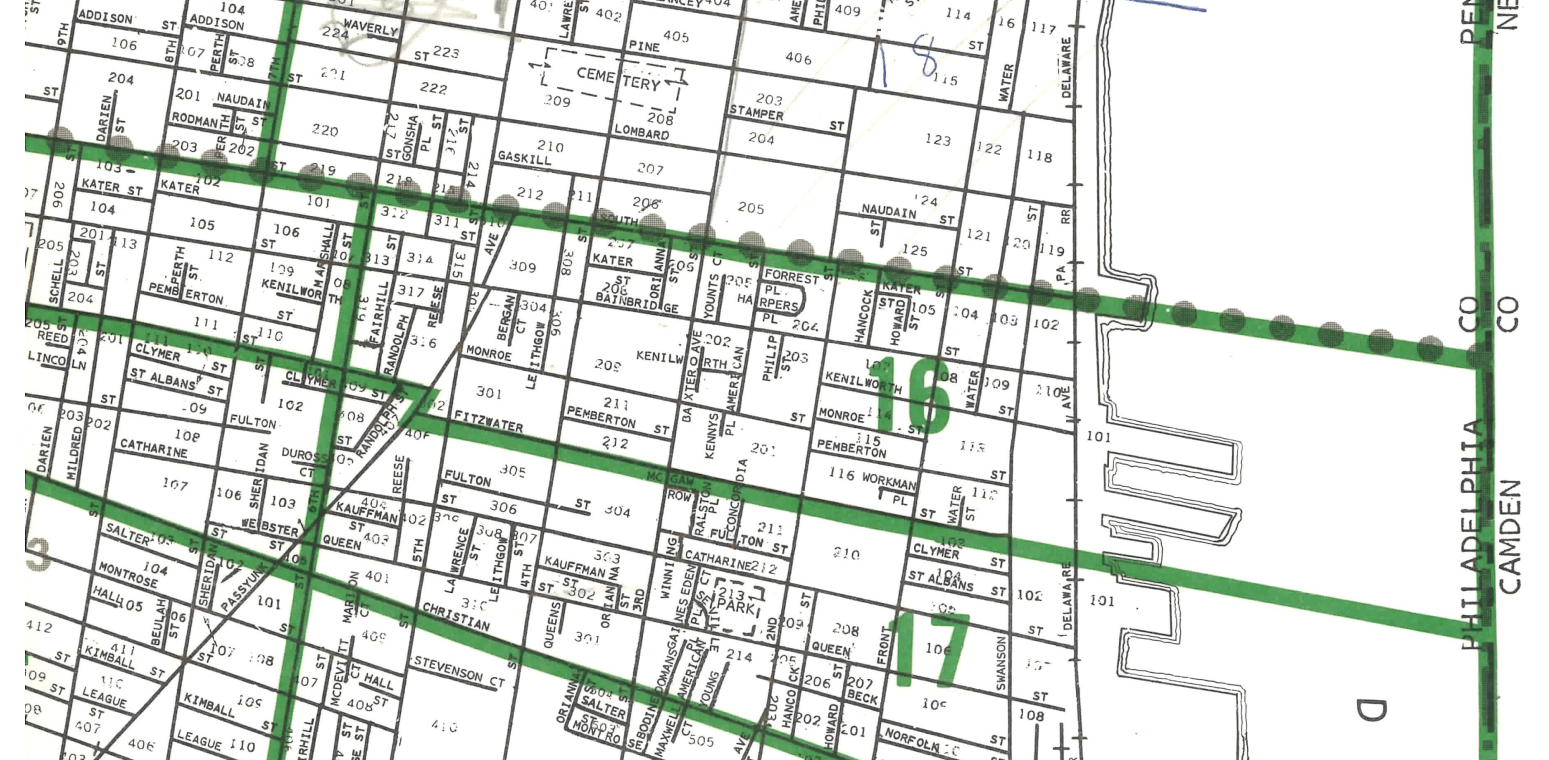
- Tract boundary segmentation is somewhat easier.

## Pitfalls

- Which block is this? Block identifiers are inconsistently located, look like street names.
- Too much detail: Block boundaries look like streets.
- "Fishhooks" are important and omnipresent.

## Our current approach

- We are relying on hand annotations for training and validating CNN.



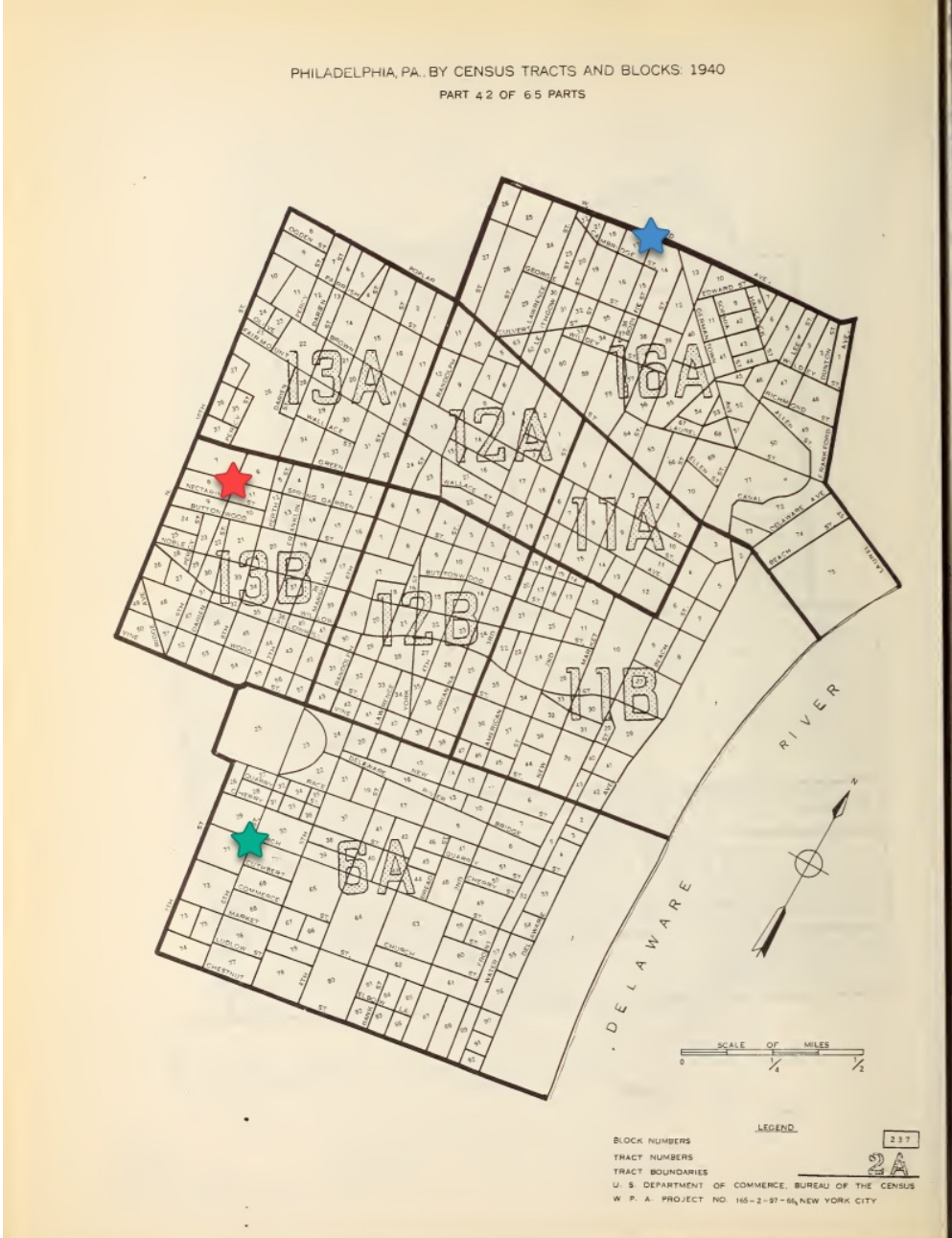
# 3 Tasks, 3 Pieces of Data



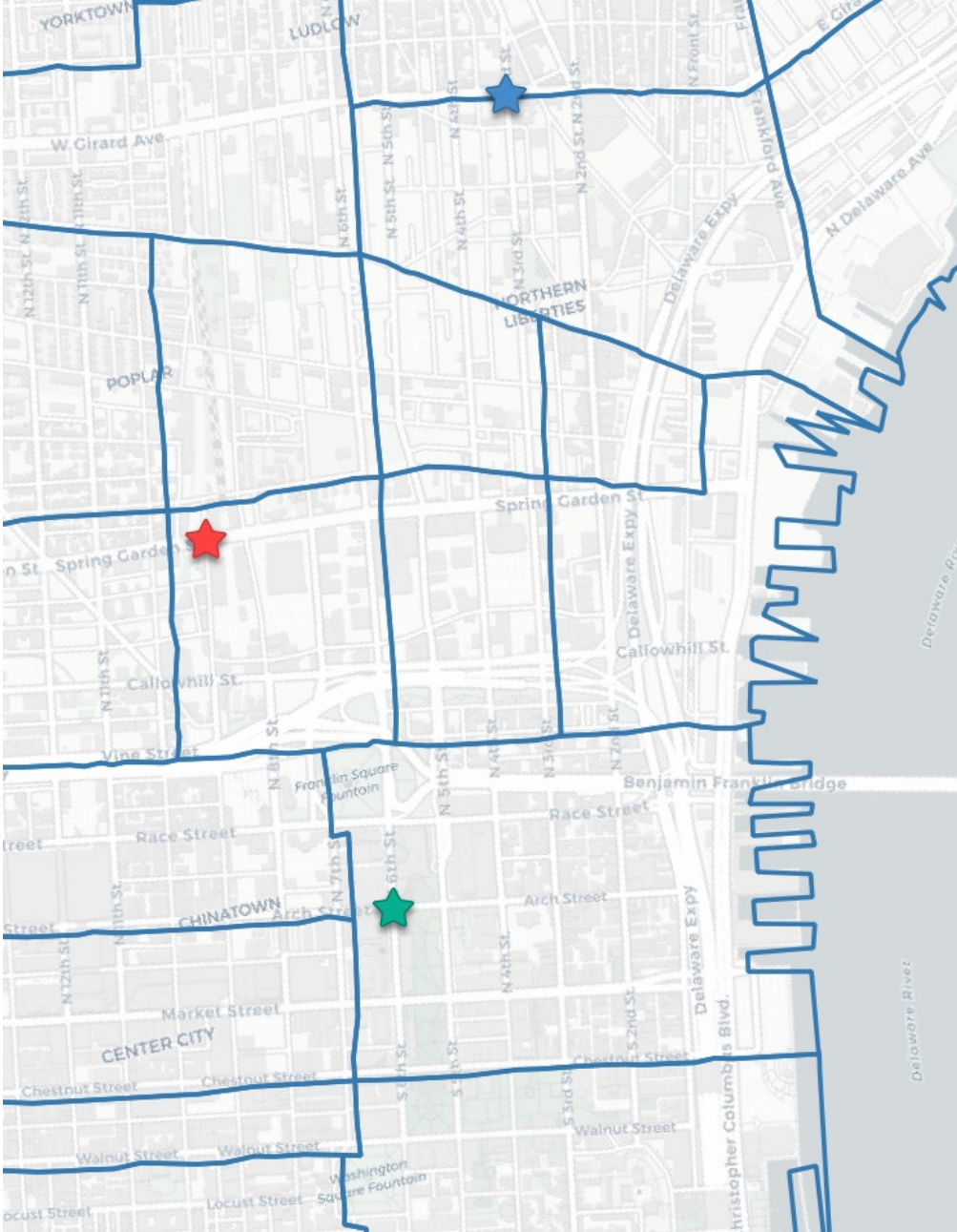
Geo-Referencing Maps



# Just Keep Clicking 🐟



+



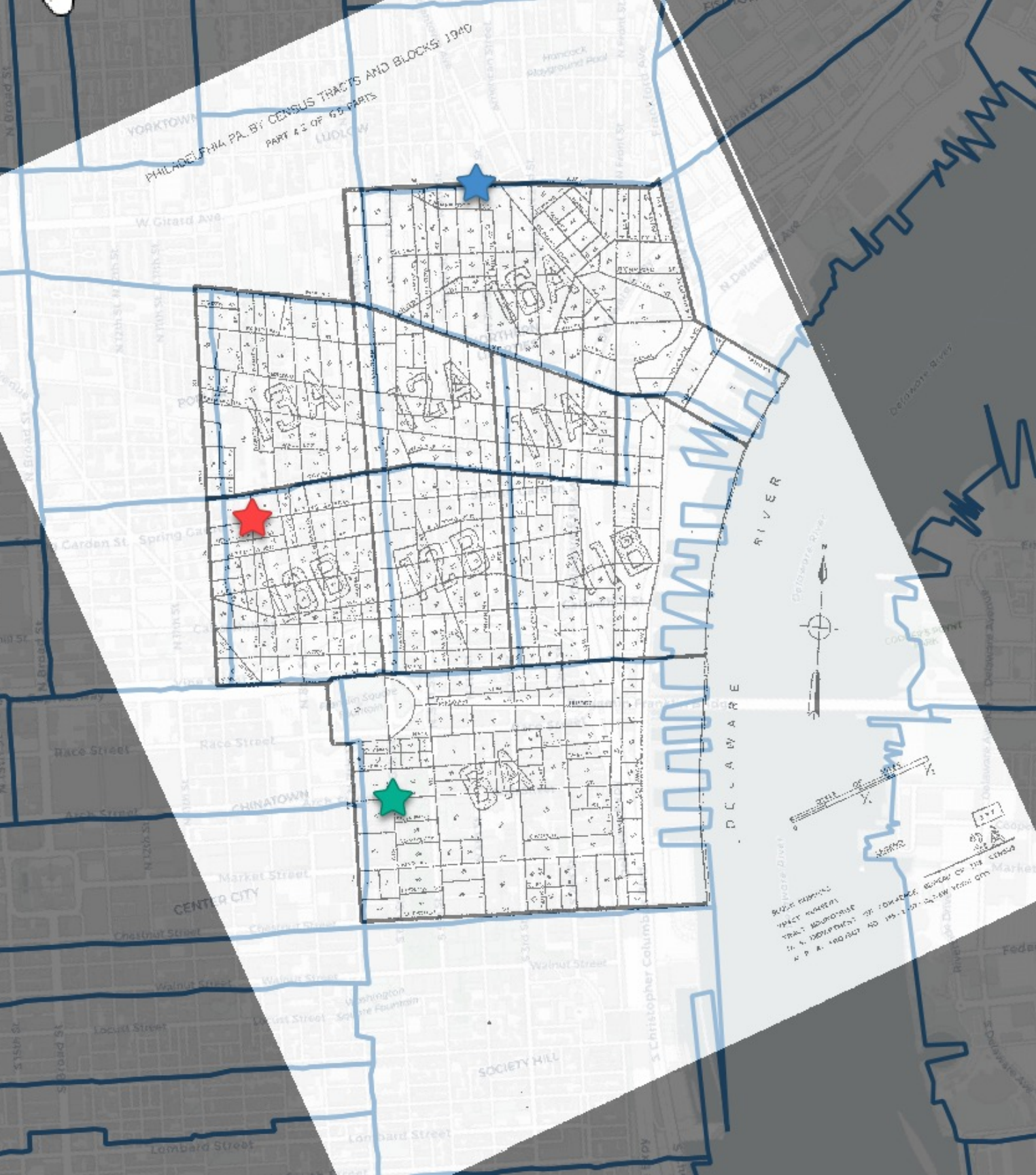
=

# A Georeferenced Map

We know where this picture goes now!

But...

- Sloooooooow.
- Only 3 points doesn't handle map inaccuracies well.
- Doesn't match (blue) reference NHGIS shapefile.



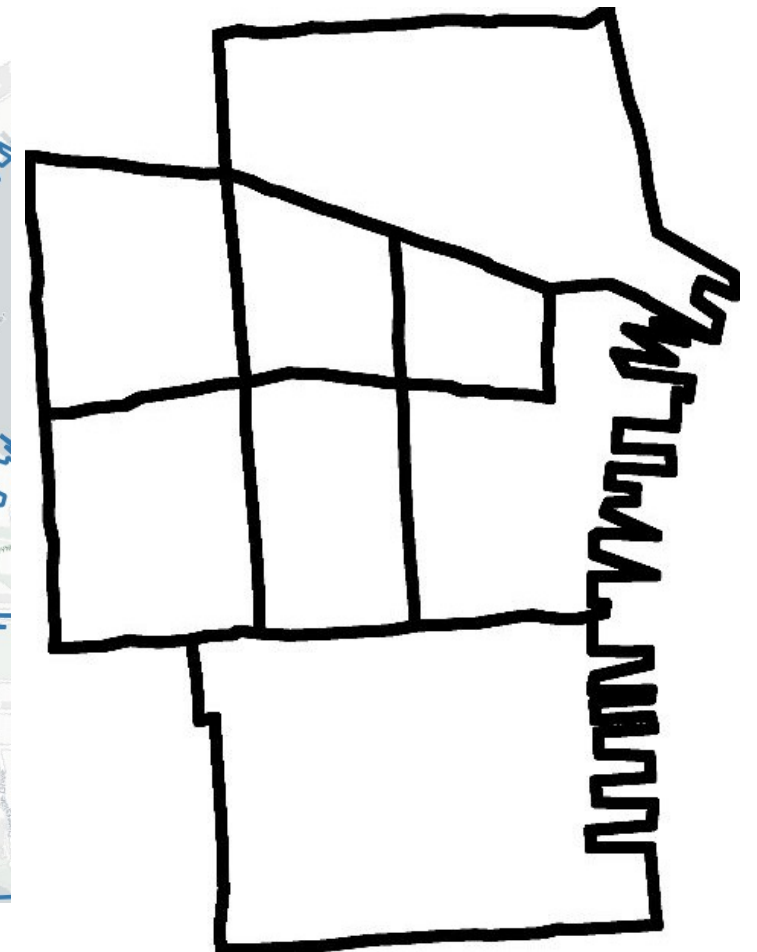
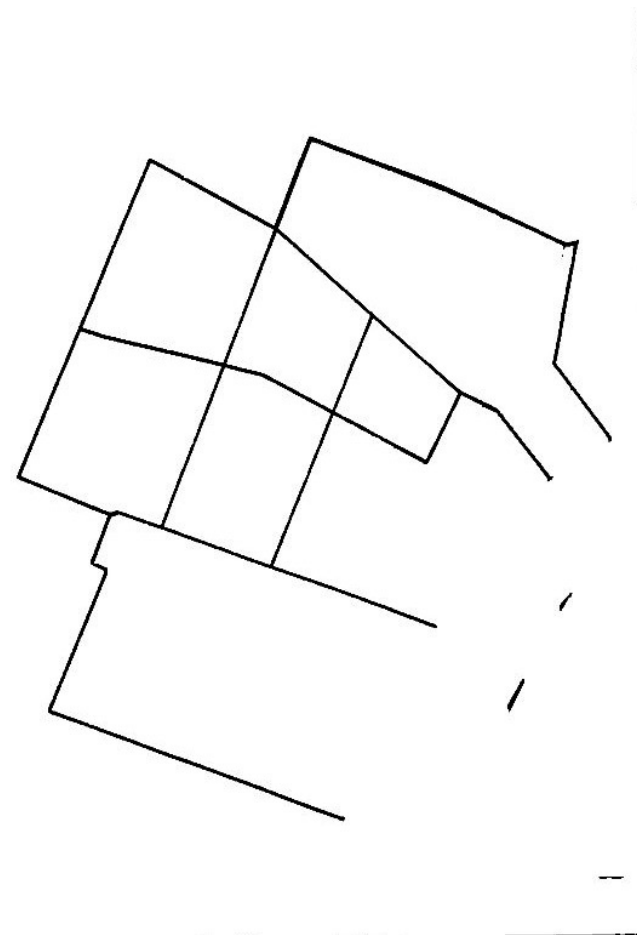
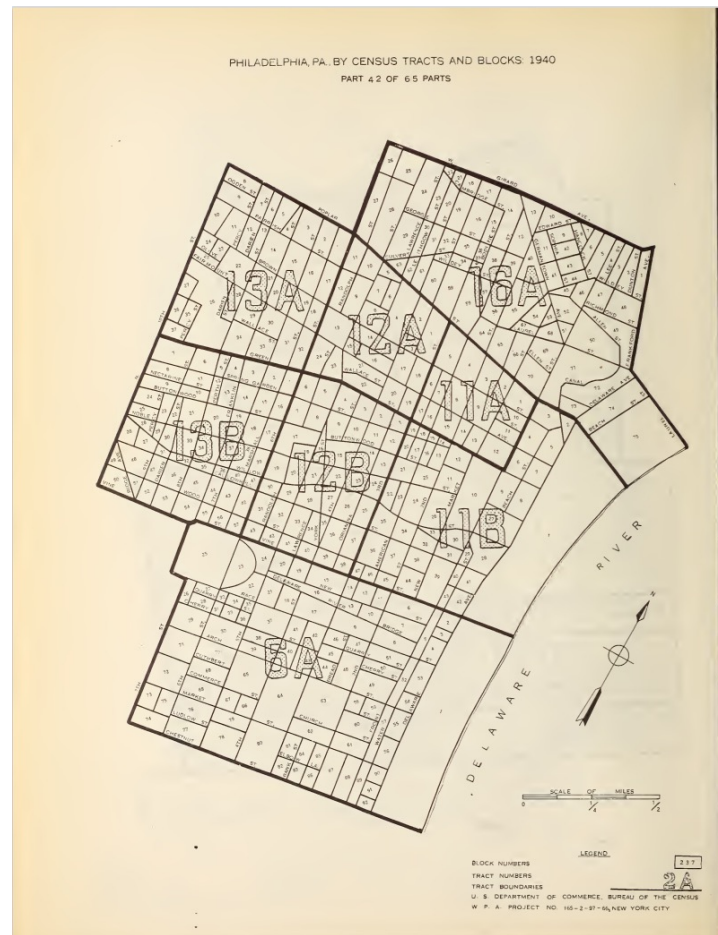
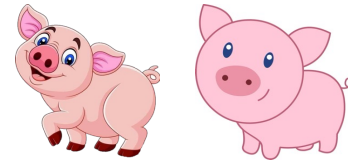




## How Can We Get Better? 🏋️‍♀️

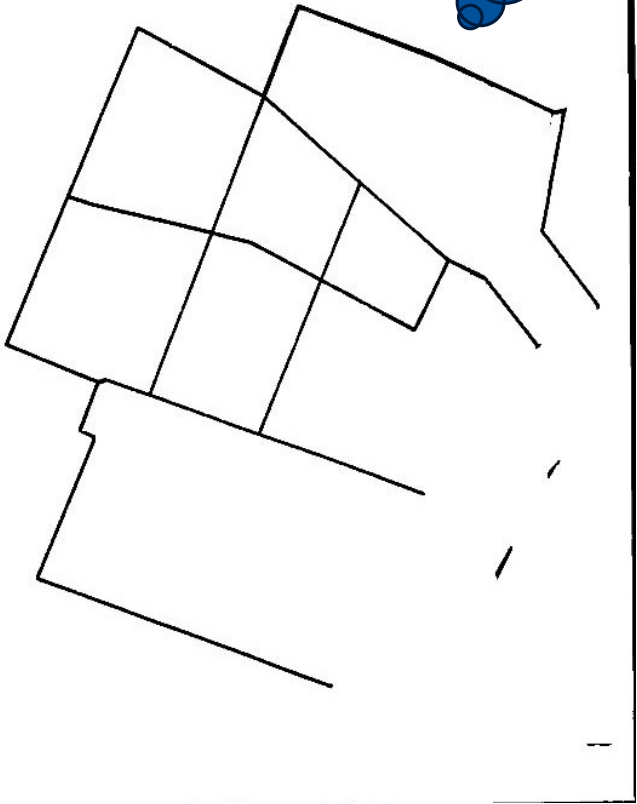
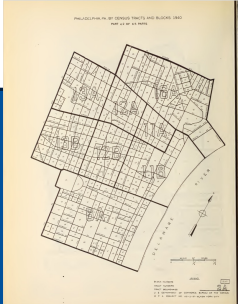
- Faster!
  - We want to process many maps.
- More accurate!
- How much of this can a computer do for us?

# Simplifying the Problem

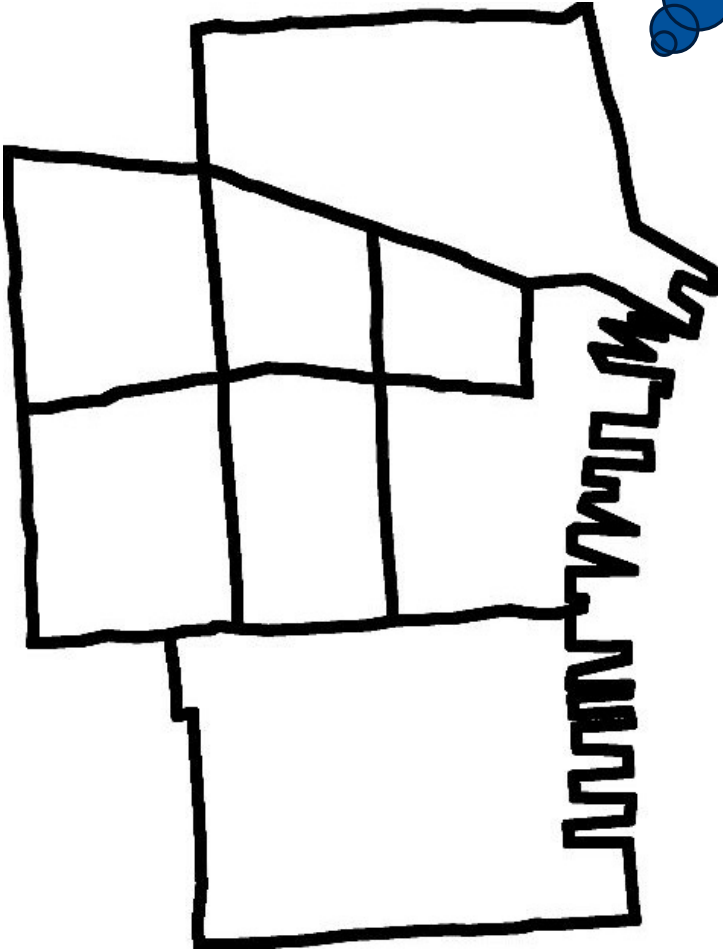


# How Hard Could This Be?

I still remember the block boundaries!

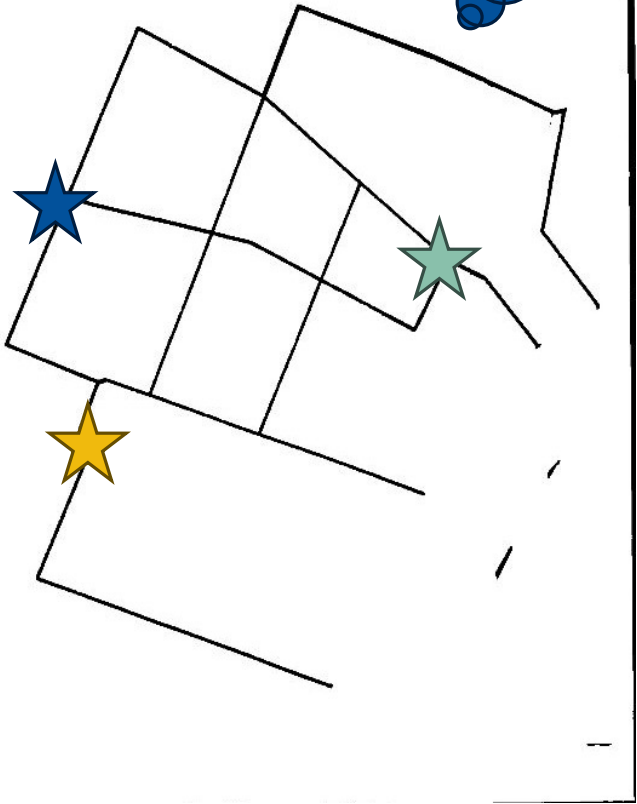
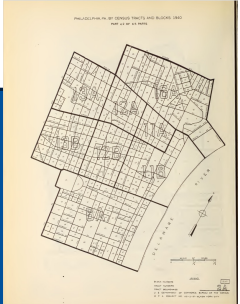


I still remember where I am in the world!

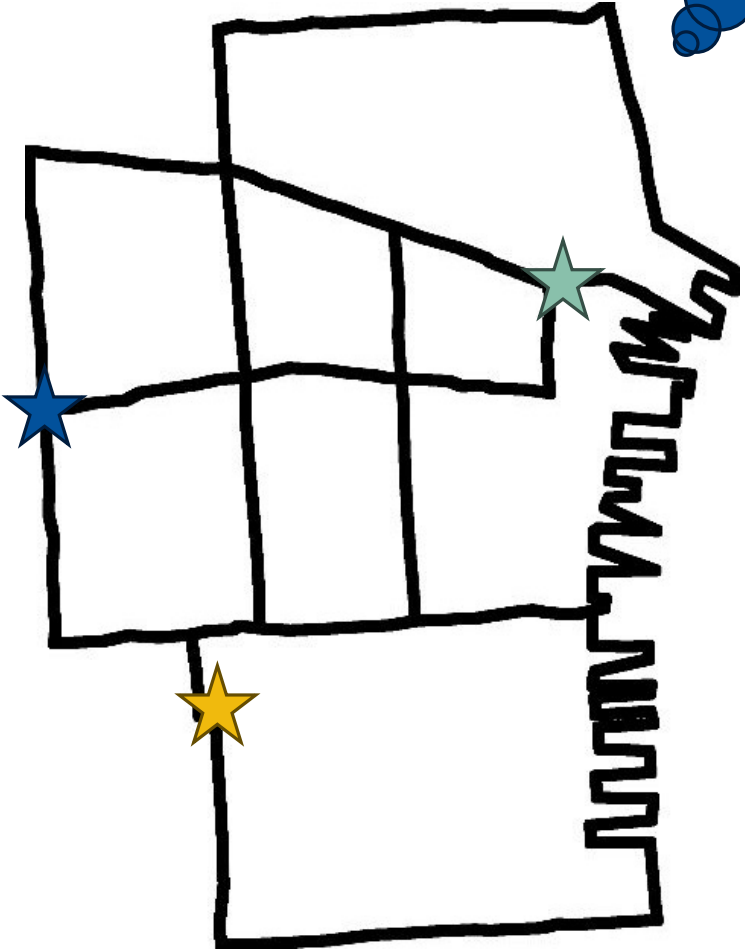


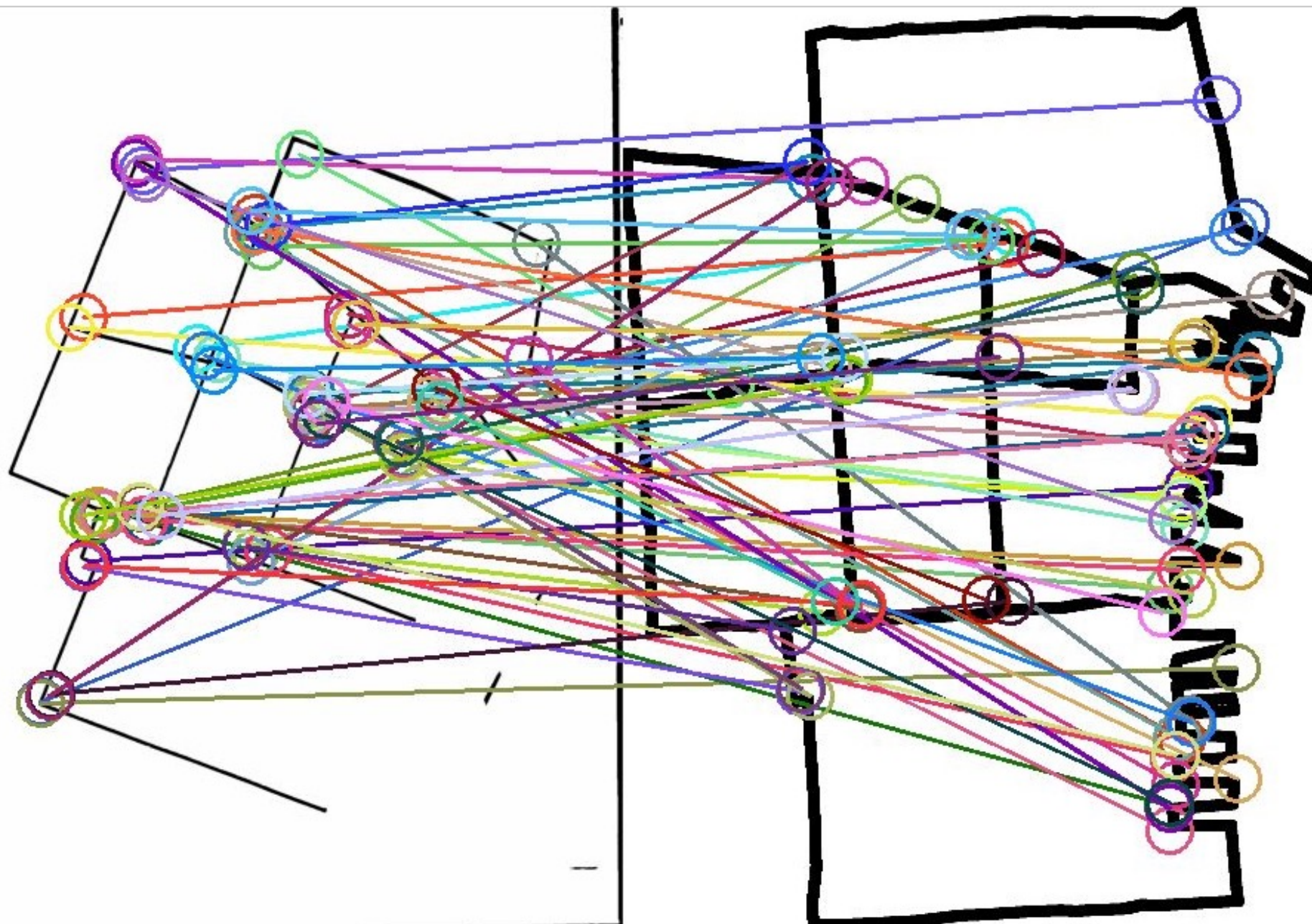
# How Hard Could This Be?

I still remember the block boundaries!



I still remember where I am in the world!





# Super. Duper. Hard.

- Too many matches.
- Several per corner!
- We need to narrow these down.

# The Basic Steps

1

Make guess

Randomly select internally consistent links.

2

Evaluate  
guess

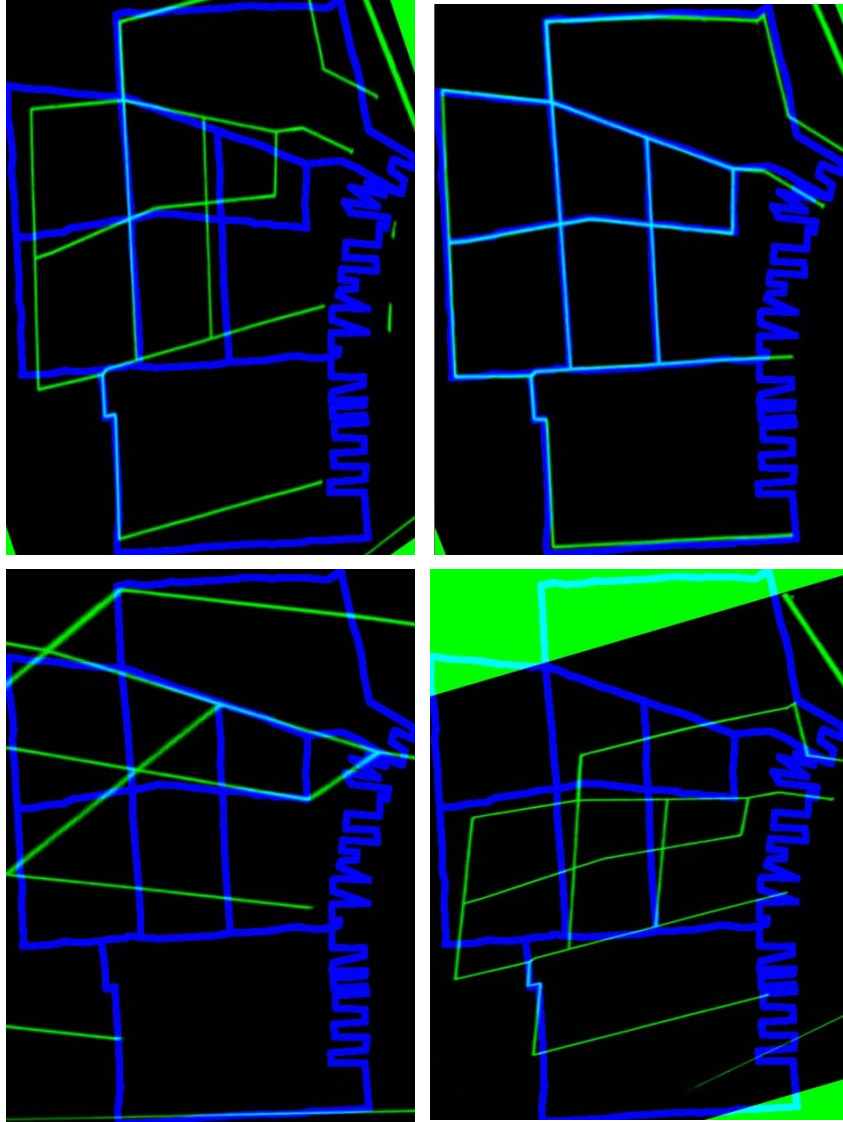
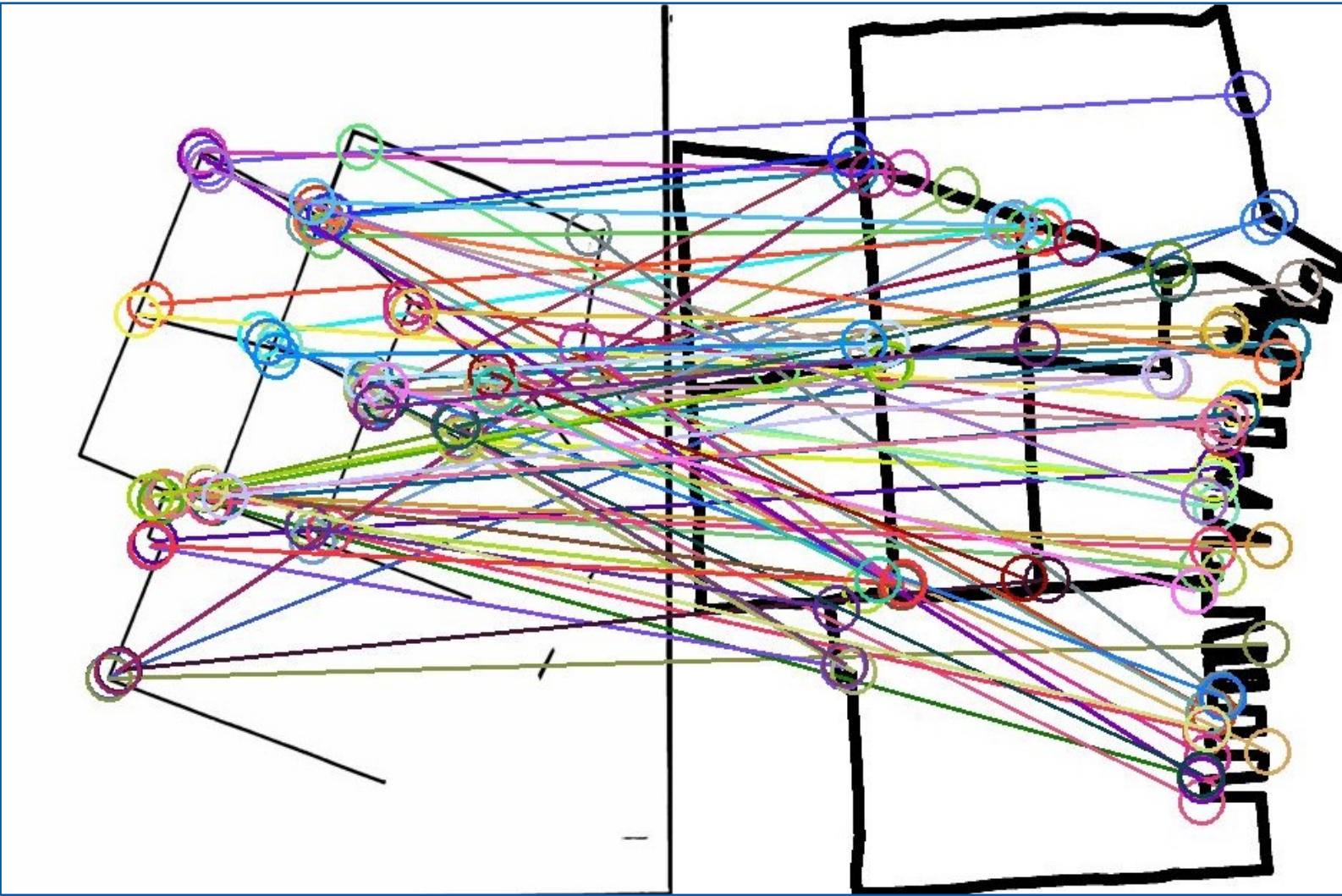
For selected links, how much do maps overlap?

3

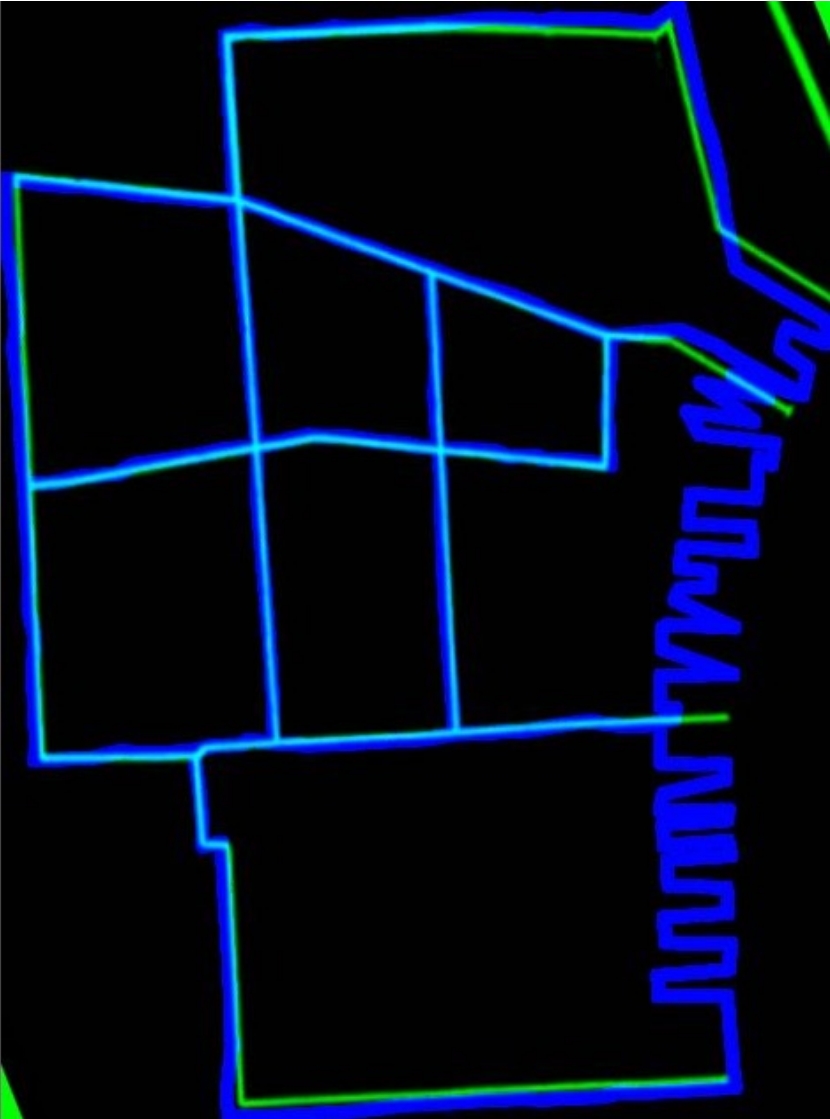
Repeat

Keep best guess.

# Selecting Links

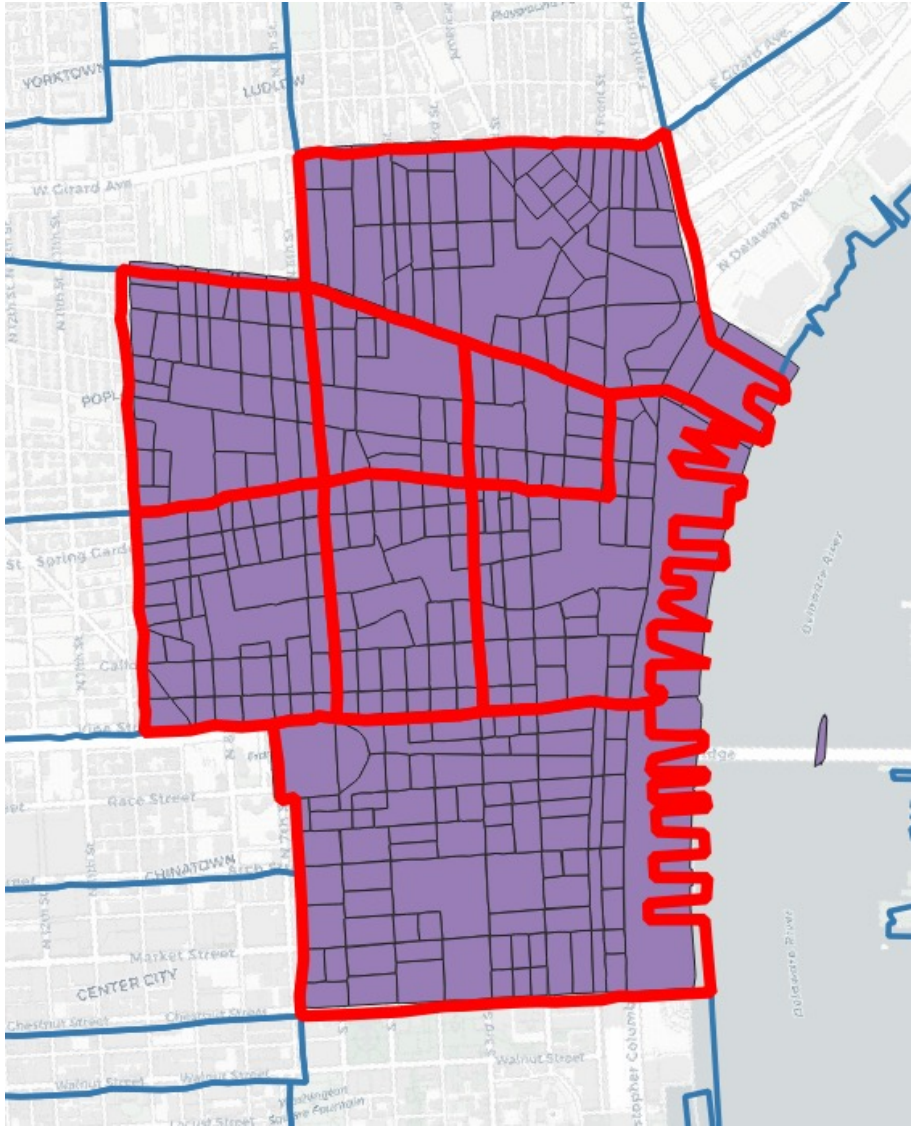
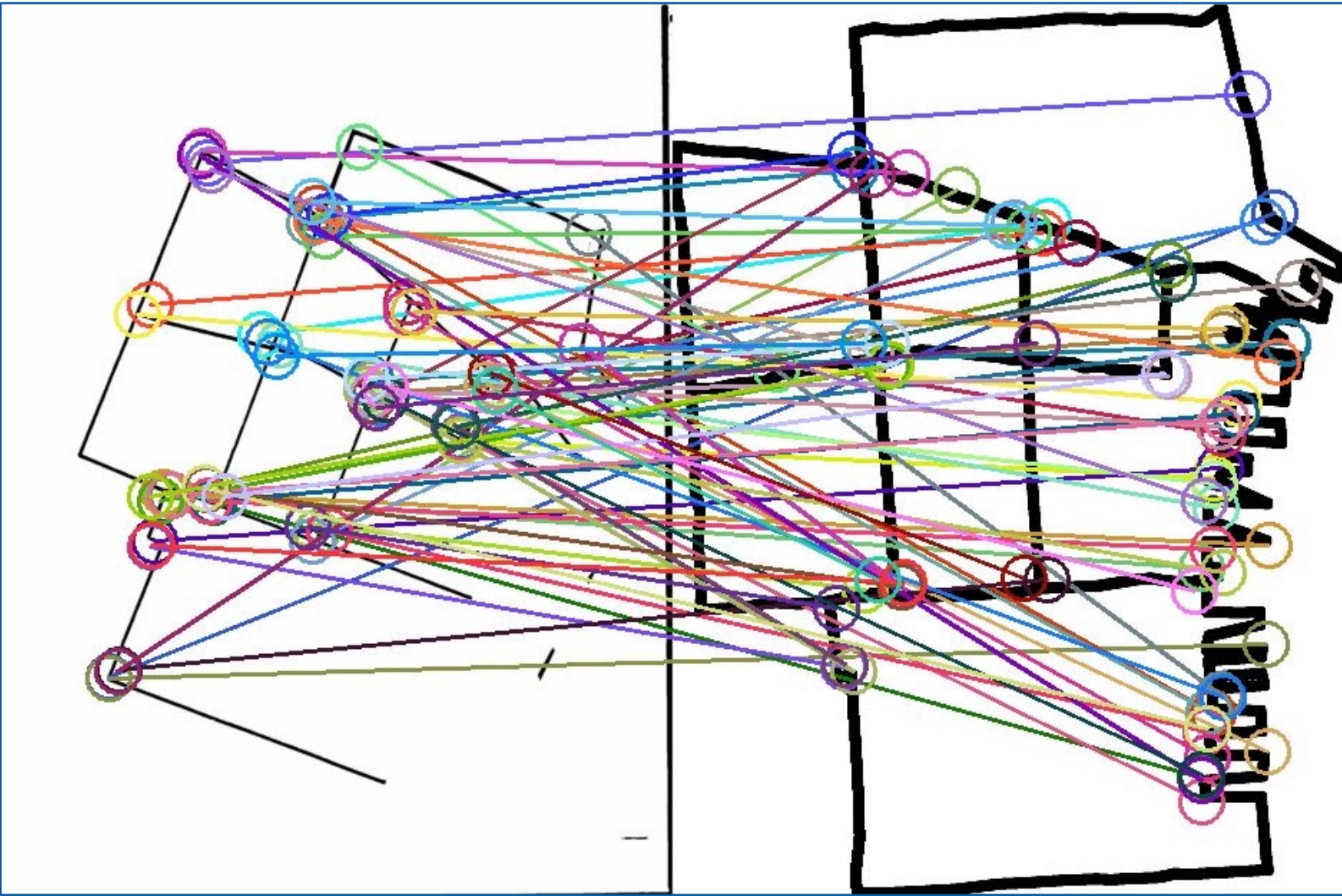


# Selecting Links





# Selecting Links



# 3 Tasks, 3 Pieces of Data

1 Shape

2 Situation

3 **Statistics**

Digitizing Tables

# Scale and Scope of Problem

- 1940, 1950, and 1960 Census of Housing, Block Statistics (1970 is digital)
- Sixteen Target Cities
  - New York City
  - Chicago
  - Philadelphia
  - Los Angeles
  - Detroit
  - Baltimore
  - Cleveland
  - St. Louis
  - Washington, DC
  - Boston
  - San Francisco
  - Pittsburgh
  - Houston
  - Cincinnati
  - Columbus, OH
  - Atlanta

# Scale and Scope of Problem

- 1940, 1950, and 1960 Census of Housing, Block Statistics
- Sixteen Target Cities
  - New York City
  - Chicago
  - Philadelphia
  - Los Angeles
  - Detroit
  - Baltimore
  - Cleveland
  - St. Louis
  - Washington, DC
  - Boston
  - San Francisco
  - Pittsburgh
  - Houston
  - Cincinnati
  - Columbus, OH
  - Atlanta
- ~2,000 pages of tabular data, ~170,000 blocks, ~2.5 million cells per decade
- Structured, tabular form, with rows and columns properly associated and with accuracy better than 99%

# Bottom Line Up Front

- Four stage process
  - Isolate table and each column
  - First pass with Tesseract
  - Algorithm to structure table
  - ML model to correct errors in OCR
- Great results
- Approach only makes sense if dataset is large and consistent

# Bottom Line Up Front (1950)

- **Custom Solution**
  - 0.07% Observations with Error
  - 0.03% Character Error Rate

# Bottom Line Up Front (1950)

- **Custom Solution**
  - 0.07% Observations with Error
  - 0.03% Character Error Rate
- **Data Entry**
  - 0.12% Observations with Error
  - 0.13% Character Error Rate
- **Tesseract (with assist with table structure)**
  - 12.94% Observations with Error
  - 7.24% Character Error Rate

# Why Not Use Out of The Box Solutions?



# Why Not Use Out of The Box Solutions?

- Adobe:

Census tract	Block	one-dwelling-structures
		Average value (dollars)
10-8	11	4.425
	12	
	15	
	16	7.288
	17	8.150
	18	
	19	
	20	
	21	
	22	
	23	
	24	6.366
	25	3.900
	26	
	28	
	29	
	30	
	31	
	32	
	34	

Cen1111 tract	Block	one-dwellInc-structure1
		Av. value (dollar)
10-e	11	4.425
11-A	12	7.288
11-e	15	8.15 0
12-A	16	6.366
	17	3,900
	18	9.50 0
	19	8,6 3 J
	20	6.483
	21	4.75 7
	22	5.87 5
	23	4.261
	24	5.34 5
	25	5,166
	26	4.80 0
	28	7.800
	29	6.0 0 0
	JO 31	4.66 6
	32	4,125

# Why Not Use Out of The Box Solutions?

- Adobe:

Census tract	Block	one-dwelling-structures
		Average value (dollars)
10-8	11	4.425
	12	
	15	
	16	7.288
	17	8.150
	18	
	19	
	20	
	21	
	22	
	23	
	24	6.366
	25	3.900
	26	
	28	
	29	
	30	
	31	
	32	
	34	

Cen1111 tract	Block	one-dwellInc-structure1
		Av. value (dollar)
10-e	11	4.425
11-A	12	7.288
11-e	15	8.15 0
12-A	16	6.366
	17	3,900
	18	9.50 0
	19	8,6 3 J
	20	6.483
	21	◀.JS 7
	22	5.87 5
	23	◀.261
	24	5.34 5
	25	5,166
	26	4.80 0
	28	7.800
	29	6.0 0 0
JO	31	◀.66 6
	32	4,125

# Why Not Use Out of The Box Solutions?

- Adobe: Bad character recognition, relation of rows lost

Census tract	Block	one-dwelling-structures Average value (dollars)	Cen1111 tract	Block	one-dwellInc- ture1 Av. value (dollar)
10-8	11	4.425	10-e	11	4.425
	12		11-A	12	7.288
	15		11-e	13	8.15 0
	16	7.288	12-A	16	6.366
	17	8.150		17	3,900
	18			18	9.50 0
	19			19	8,6 3 J
	20			20	6.483
	21			21	◀.JS 7
	22			22	5.87 5
	23			23	◀.261
	24	6.366		24	5.34 5
	25	3.900		25	5,166
	26			26	4.80 0
	28			28	7.800
	29			29	6.0 0 0
	30			JO 31	◀.66 6
	31			32	4,125
	32				
	34				

# Why Not Use Out of The Box Solutions?

- Textract:

Census tract	Block	All dwelling units by occupancy and tenure					All dwelling units by condition and plumbing facilities			Occupied dwelling units			Contract monthly rent <sup>1</sup>		Value <sup>2</sup> of one-dwelling-unit structures		
		Total	Owner occupied	Renter occupied	Vacant non-seasonal not dilap., for rent or sale	Other vacant and non-resident	Number reporting	No private bath or dilap.	No running water or dilap.	Total	Persons per room		Occupied by non-white	Number reporting	Average monthly rent (dollars)	Number reporting	Average value (dollars)
											Number reporting	1.51 or more					
46-G	15	91	42	46	2	1	75			88	86	3		41	44.34	23	8,630
	16	80	44	34		2	71	2	2	78	77	4		34	43.20	32	16,050
	17	100	36	57	5	2	94	1	1	93	93			55	48.83	15	10,866
	18	75	43	30		2	66	1	1	73	73	4		30	50.93	31	10,064
	19	79	66	11	2		72	1	1	77	77	2		7	38.28	40	7,040
	20	75	42	19	4	10	56			61	56	1		16	44.43	32	6,921
	21	66	59	7			62	3	2	66	66	1		7	44.85	48	7,479
	22	42	35	7			42	1	1	42	42			6	42.00	30	8,300
	23	50	46	1		3	42			47	44					26	8,192
	24	64	53	10	1		62	6	6	63	63			10	53.10	23	7,434
	25	57	47	10			54			57	55	1		9	40.22	40	9,575
	26	106	47	58		1	98	2	2	105	97	2		57	39.22	20	9,600
	27	64	39	25			58	3	3	64	64	1		24	41.29	30	6,543
	28	77	32	44		1	76	1	1	76	76	2		43	43.88	20	7,570
	29	79	39	36	3	1	59	2	2	75	69	3		33	46.12	23	7,869
	30	58	42	10	1	5	45			52	47	1		9	47.33	26	8,323
	31	49	42	7			49	1		49	49	1		6	46.83	41	8,402
32	62	45	17			60	5		62	61	2		16	46.43	39	8,128	

# Why Not Use Out of The Box Solutions?

- Textract: when it works it works!
- 1.5% error rate, 0.22% ignoring cell alignment errors (stats for this page only)

Census tract	Block	All dwelling units by occupancy and tenure					All dwelling units by condition and plumbing facilities			Occupied dwelling units			Contract monthly rent		Value of one-dwelling-unit structures	
		Total	Owner occupied	Renter occupied	Vacant non-seasonal not dilap. for rent or sale	Other vacant and non-resident	Number reporting	No private bath or dilap.	No running water or dilap.	Total	Persons per room	Occupied by non-white	Number reporting	Average monthly rent (dollars)	Number reporting	Average value (dollars)
46-G	15	91	42	46	2	1	75			88	86	3	41	4434	23	8630
	16	86	44	34		2	71	2	2	76	77	4	34	4320	32	16050
	17	100	36	57	5	2	94	1	1	93	93		55	4883	15	10866
	18	75	43	30		2	66	1	1	73	73	4	30	5093	31	10064
	19	79	66	11	2		72	1	1	77	77	2	7	3828	40	7040
	20	75	42	19	4	10	56			61	56	1	16	4443	32	6921
	21	66	59	7			62	3	2	66	66	1	7	4485	48	7479
	22	42	33	7			42	1	1	42	42		6	4200	30	8300
	23	50	46	1		3	42			47	44				26	8192
	24	64	53	10	1		62	6	6	63	63		10	5310	23	7434
	25	57	47	10			54			57	55	1	9	4022	40	9575
	26	106	47	58		1	98	2	2	105	97	2	57	3922	20	9600
	27	64	39	25			56	3	3	64	64	1	24	4129	30	6543
	28	77	32	44		1	76	1	1	76	76	2	43	4386	20	7576
	29	79	39	36	3	1	59	2	2	75	69	3	33	4612	23	7869
	30	58	42	10	1	5	45			52	47	1	9	4733	26	8323
	31	49	42	7			49	1		49	49	1	6	4683	41	8402
	32	62	45	17			60	5		62	61	2	16	4643	39	8128

# Why Not Use Out of The Box Solutions?

- Textract: when it doesn't work...

Census tract	Block	All dwelling units by occupancy and tenure					All dwelling units by condition and plumbing facilities			Occupied dwelling units			Contract monthly rent <sup>1</sup>		Value <sup>2</sup> of one-dwelling-unit structures		
		Total	Owner occupied	Renter occupied	Vacant non-seasonal not dilap., for rent or sale	Other vacant and non-resident	Number reporting	No private bath or dilap.	No running water or dilap.	Total	Persons per room		Occupied by non-white	Number reporting	Average monthly rent (dollars)	Number reporting	Average value (dollars)
											Number reporting	1.51 or more					
33-1	2	60	43	17					60	60	3		17	32.41	42	4,821	
	5	35	27	8					35	34	1		8	34.62	19	4,410	
	6	77	55	22					77	75			22	35.45	46	5,260	
	7	37	30	6	1				37	36			7	31.71	24	5,020	
	8	24	22	2					24	24			2		17	5,500	
	9	26	21	5					26	26			5	31.00	19	5,894	
	11	10	7	3					10	10			3	61.66	7	6,142	
	12	49	7	40	1	1			47	46	4		41	25.70	5	3,180	
	14	9	6	3					9	8			3	40.00	6	7,483	
	15	49	32	17					49	49	1		17	26.35	32	3,728	
	16	56	34	20	1	1			56	54	1		20	38.30	25	4,980	
	17	18	7	11					18	18			11	35.18	4	4,875	
	18	19	11	8					19	19			8	30.12	11	6,045	
	19	38	34	4					38	37			4	30.50	32	5,812	
	20	30	23	7					30	30			4	33.00	22	7,463	
	21	63	40	22		1			62	59	1		20	29.05	35	5,200	
	22	45	17	28					45	44			28	29.48	16	7,500	
	23	26	25	1					26	26			1		23	7,026	
	24	26	19	6	1				26	25			6	37.50	19	6,868	

# Why Not Use Out of The Box Solutions?

- Textract: when it doesn't work... it doesn't work! Small input tweaks do not fix error.

Census tract	Block	All dwelling units by occupancy and tenure				All dwelling units by condition and plumbing facilities			Occupied dwelling units			Contract monthly rent		Value of one-dwelling-unit structures	
		Total	Owner occupied	Renter occupied	Vacant non-seasonal not for rent or sale	Number reporting	No private bath or dilap.	No running water or dilap.	Total	Persons per room	Number reporting	Average monthly rent (dollars)	Number reporting	Average value (dollars)	
38-1	1	60	35	25	0	60	0	0	60	1.51	17	3241	42	4821	
	2	37	22	15	0	34	0	0	34	1.51	18	3462	19	4410	
	3	77	55	22	0	76	0	0	77	1.51	22	3545	46	5260	
	4	37	30	7	0	37	0	0	36	1.51	7	2121	24	5020	
	5	24	22	2	0	24	0	0	24	1.51	25	3100	19	5500	
	6	26	21	5	0	26	0	0	26	1.51	3	6166	7	5894	
	7	10	7	3	0	10	0	0	10	1.51	3	2570	5	6142	
	8	41	32	9	0	41	0	0	47	1.51	4	1048	1	3180	
	9	9	3	6	0	9	0	0	9	1.51	3	1048	1	3180	
	10	56	47	9	0	56	0	0	54	1.51	2	3830	25	4980	
	11	18	11	7	0	18	0	0	18	1.51	1	3518	4	4875	
	12	19	11	8	0	19	0	0	19	1.51	8	3012	11	6045	
	13	38	33	5	0	38	0	0	37	1.51	4	3050	4	5812	
	14	30	23	7	0	30	0	0	30	1.51	4	3300	22	7463	
	15	60	40	20	0	60	0	0	62	1.51	20	2905	35	5200	
	16	45	28	17	0	45	0	0	44	1.51	28	2948	16	7500	
	17	26	17	9	0	26	0	0	26	1.51	1	2948	1	7500	

# Why Not Use Out of The Box Solutions?

- Textract:
  - Sample Size: 169 Pages
  - Catastrophic Failures: 45
  - Moderate Failures: 6
  - **Unacceptable *page level* error: 30%**
  - Small errors in table layout can be algorithmically corrected, catastrophic failures cannot



# Method

- Isolate table
- Isolate columns
- Tesseract columns
- Structure into table
- Match to labeled data
- Train model to correct Tesseract errors
- Visualize and correct issues throughout
- Final check for internal consistency and vs tract

# Isolate Table Body, Straighten Image

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

["Total population" contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more]

Blocks within census tracts	Total population	All housing units by condition and plumbing									Occupied housing units							
		Total	Sound			Deteriorating			Dilapidated	Owner occupied			Renter occupied			Occu- pied by non- white	1.01 or more per- sons per room	
			Total	With all plumb- ing facil- ities	Lack- ing some or all facil- ities	Total	With all plumb- ing facil- ities	Lacking some or all facilities		With flush toilet	No flush toilet	Total	Average value (dollars)	Average num- ber of rooms	Total			Average con- tract rent (dollars)
21...	30	9	8	8	...	1	1	...	...	...	6	3500	5.3	3	...	...	...	...
22...	87	28	22	22	...	6	1	5	...	...	14	6000	7.1	12	43	4.2	...	...
23...	124	44	34	30	4	9	4	...	5	1	22	6000	7.0	16	28	5.3	...	3
24...	**120	35	26	24	2	9	9	...	...	...	20	4500	7.1	12	40	4.4	...	1
25...	247	58	38	38	...	15	14	1	...	2	41	5000	7.5	11	47	5.6	...	2
26...	145	40	25	25	...	10	9	1	...	5	28	5000	6.9	8	39	5.9	...	5
27...	85	20	12	12	...	7	7	...	...	1	16	5000	7.0	4	...	...	...	5
28...	21	7	2	2	...	5	5	...	...	...	1	...	...	4	...	...	...	2
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	2
30...	51	13	12	12	...	...	...	...	...	1	9	4000	5.9	2	...	...	...	2
31...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
34...	10	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
35...	9	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36...	54	18	7	7	...	8	8	...	...	3	5	6500	8.4	4	...	...	...	...
37...	60	15	11	11	...	4	3	1	...	...	9	5500	7.6	6	37	6.2	...	2
38...	133	52	38	36	2	13	8	5	...	1	31	5000	7.4	9	41	5.3	...	3
43...	35	12	5	5	...	4	2	1	1	3	2	...	...	5	29	6.0	...	2
44...	6	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	3
45...	3	1	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
49...	141	43	17	16	1	26	20	3	3	...	15	5500	8.3	21	31	4.9	1	4
50...	84	35	14	13	1	21	4	10	7	...	12	3500	5.0	15	22	4.4	...	4
51...	22	9	7	7	...	2	1	1	...	...	5	...	7.6	1	...	...	...	...
18-B....	*8802	2781	2358	2287	71	338	248	57	33	85	1823	5000	6.5	727	40	4.7	2	171
1....	111	53	38	33	5	14	12	2	...	1	18	5000	6.6	28	...	...	...	...

# Isolate Table Body, Straighten Image

Pass through Textract

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

["Total population" contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more]

Blocks within census tracts	Total population	All housing units by condition and plumbing							Occupied housing units							
		Total	Sound		Deteriorating			Dilapidated	Owner occupied		Renter occupied			Occu- pied by non- white	1.01 or more per- sons per room	
			With all plumb- ing facil- ities	Lack- ing some or all facil- ities	Total	With all plumb- ing facil- ities	Lacking some or all facilities		No flush toilet	Total	Average value (dollars)	Average number of rooms	Total			Average con- tract rent (dollars)
21...	30	9	8	...	1	11	...	...	6	3500	5.3	3	...	...	...	
22...	87	28	22	20	6	11	...	...	22	6000	7.1	12	...	...	...	
23...	124	44	34	30	9	11	...	...	14	6000	7.0	16	...	...	...	
24...	**120	35	26	24	2	9	...	...	20	4500	7.1	12	...	...	...	
25...	247	55	38	38	...	15	...	...	41	5000	7.5	11	...	...	...	
26...	145	40	28	25	...	10	...	...	28	5000	6.9	8	...	...	...	
27...	85	20	12	12	...	7	...	...	16	5000	7.0	4	...	...	...	
28...	21	7	2	2	...	5	...	...	1	...	...	4	...	...	...	
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...	
30...	51	13	12	12	...	...	...	...	9	4000	5.9	2	...	...	...	
31...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...	
34...	10	2	...	...	...	...	...	...	...	...	...	...	...	...	...	
35...	9	2	...	...	...	...	...	...	...	...	...	...	...	...	...	
36...	54	18	7	7	...	8	...	...	5	6500	8.4	4	...	...	...	
37...	60	15	11	11	...	4	...	...	9	5500	7.6	6	...	...	...	
38...	133	52	38	36	2	13	...	...	31	5000	7.4	9	...	...	...	
43...	35	12	5	5	...	4	...	...	2	...	...	5	...	...	...	
44...	6	2	...	...	...	...	...	...	...	...	...	...	...	...	...	
45...	3	1	...	...	...	...	...	...	...	...	...	...	...	...	...	
49...	141	43	17	16	1	26	...	...	15	5500	8.3	21	...	...	...	
50...	84	35	14	13	1	21	...	...	12	3500	5.0	15	...	...	...	
51...	22	6	7	7	...	2	...	...	5	...	7.6	1	...	...	...	
18-B....	*8802	2781	2358	2287	71	338	248	57	85	1823	5000	6.5	727	40	4.7	171
1....	111	53	38	38	...	14	12	2	1	18	5000	6.6	25	...	...	...

# Isolate Table Body, Straighten Image

Find where, vertically, we go from mostly alpha to mostly numeric

Mostly Text

Mostly Numbers

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

["Total population" contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more]

Blocks within census tracts	Total population	All housing units by condition and plumbing							Occupied housing units							
		Total	Sound		Deteriorating			Dilapidated	Owner occupied			Renter occupied			Occ. by non-white	1.01 or more persons per room
			With all plumbing facilities	Lacking some or all facilities	Total	With all plumbing facilities	Lacking some or all facilities		With flush toilet	No flush toilet	Total	Average value (dollars)	Average number of rooms	Total		
21...	30	9	8	...	1	1	...	...	6	3500	5.3	3	...	...	...	...
22...	87	28	22	22	8	...	1	...	22	6000	7.1	...	...	...	...	...
23...	124	44	34	30	4	9	4	...	14	6000	7.0	12	43	4.2	...	...
24...	**120	35	26	24	2	9	9	...	20	4500	7.1	16	28	5.3	...	...
25...	247	55	38	38	...	15	14	...	41	5000	7.5	12	40	4.4	...	...
26...	145	40	25	25	...	10	9	...	28	5000	6.9	1	47	5.6	...	...
27...	85	20	12	12	...	7	7	...	16	5000	7.0	8	39	5.9	...	...
28...	21	7	2	2	...	5	5	...	1	...	...	4	...	...	...	...
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...
30...	51	13	12	12	...	...	...	...	9	4000	5.9	...	...	...	...	...
31...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...
34...	10	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...
35...	9	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36...	54	18	7	7	...	8	8	...	5	6500	8.4	4	...	...	...	...
37...	60	15	11	11	...	4	3	...	9	5500	7.6	6	37	6.2	...	...
38...	133	52	38	36	2	13	8	...	31	5000	7.4	9	41	5.3	...	...
43...	35	12	5	5	...	4	2	...	2	...	...	5	29	6.0	...	...
44...	6	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
45...	3	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
49...	141	43	17	16	1	26	20	...	15	5500	8.3	21	31	4.9	...	...
50...	84	35	14	13	1	21	4	...	12	3500	5.0	15	22	4.4	...	...
51...	22	6	7	7	...	2	1	...	5	...	7.6	1	...	...	...	...
18-B....	*8802	2781	2358	2287	71	338	248	57	1823	5000	6.5	727	40	4.7	2	171
1....	111	53	38	38	...	14	12	...	18	5000	6.6	...	...	...	...	...

# Isolate Table Body, Straighten Image

Table is isolated

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

“Total population” contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more.

Blocks within census tracts	Total population	All housing units by condition and plumbing							Occupied housing units						
		Total	Sound		Deteriorating		Dilapidated	Owner occupied		Renter occupied			Occupied by non-white	1.01 or more persons per room	
			Total	With all plumbing facilities	Lacking some or all facilities	Total		With all plumbing facilities	Lacking some or all facilities	Total	Average value (dollars)	Average number of rooms			Average contract rent (dollars)
21...	30	9	8	8	...	11	11	...	6	3500	5.3	3	...	...	...
22...	87	28	22	22	...	11	11	...	14	6000	7.1	12	...	...	...
23...	124	44	34	30	4	9	4	...	22	6000	7.0	16	...	...	...
24...	**120	35	26	24	2	9	9	...	20	4500	7.1	12	...	...	...
25...	247	53	38	38	...	15	14	...	41	5000	7.5	11	...	...	...
26...	145	40	25	25	...	10	9	...	28	5000	6.9	8	...	...	...
27...	85	20	12	12	...	7	7	...	16	5000	7.0	4	...	...	...
28...	21	7	2	2	...	5	5	...	1	...	...	4	...	...	...
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...
30...	51	13	12	12	...	...	...	...	9	4000	5.9	2	...	...	...
31...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...
34...	10	...	...	...	...	...	...	...	...	...	...	...	...	...	...
35...	9	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36...	54	18	7	7	...	8	8	...	5	6500	8.4	4	...	...	...
37...	60	15	11	11	...	4	3	...	9	5500	7.6	6	...	...	...
38...	133	52	38	36	2	13	8	...	31	5000	7.4	9	...	...	...
43...	35	12	5	5	...	4	2	...	2	...	...	5	...	...	...
44...	6	...	...	...	...	...	...	...	...	...	...	...	...	...	...
45...	3	...	...	...	...	...	...	...	...	...	...	...	...	...	...
49...	141	43	17	16	1	26	20	...	15	5500	8.3	21	...	...	...
50...	84	35	14	13	1	21	4	...	12	3500	5.0	15	...	...	...
51...	22	6	7	7	...	2	1	...	5	...	7.6	1	...	...	...
18-B....	*8802	2781	2358	2287	71	338	248	57	85	1823	5000	6.5	727	40	4.7
1....	111	53	58	33	5	14	12	2	11	18	5000	6.6	28	17	1.7

# Isolate Table Body, Straighten Image

Find the rotated bounding box that contains all the body bounding boxes

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

["Total population" contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more]

Blocks within census tracts	Total population	All housing units by condition and plumbing							Occupied housing units								
		Total	Sound		Deteriorating			Dilapidated	Owner occupied			Renter occupied			Occu- pied by non- white	1.01 or more per sons per room	
			With all plumb- ing facil- ities	Lack- ing some or all facil- ities	Total	With all plumb- ing facil- ities	Lacking some or all facilities		With flush toilet	No flush toilet	Total	Average value (dollars)	Aver- age num- ber of rooms	Total			Average con- tract rent (dollars)
21...	30	9	8	8	...	1	1	...	...	6	3500	5.3	3	...	...	...	...
22...	87	28	22	22	...	2	1	...	...	14	6000	7.1	12	...	...	...	...
23...	124	44	34	30	4	9	4	...	1	22	6000	7.0	16	...	...	...	...
24...	**120	35	26	24	2	9	9	...	...	20	4500	7.1	12	...	...	...	...
25...	247	55	38	38	...	15	14	...	2	41	5000	7.5	1	...	...	...	...
26...	145	40	28	25	...	10	9	...	1	28	5000	6.9	8	...	...	...	...
27...	85	20	12	12	...	7	7	...	1	16	5000	7.0	4	...	...	...	...
28...	21	7	2	2	...	5	5	...	...	1	...	...	4	...	...	...	...
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
30...	51	13	12	12	...	...	...	...	1	9	4000	5.9	2	...	...	...	...
31...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
32...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
33...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
34...	10	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
35...	9	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36...	54	18	7	7	...	8	8	...	3	5	5500	8.4	4	...	...	...	...
37...	60	15	11	11	...	4	3	...	...	9	5500	7.6	6	...	...	...	...
38...	133	52	38	36	2	13	8	...	1	31	5000	7.4	9	...	...	...	...
43...	35	12	5	5	...	4	2	...	1	2	...	...	5	...	...	...	...
44...	6	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
45...	3	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
49...	141	43	17	16	1	26	20	...	3	15	5500	8.3	21	...	...	...	...
50...	84	35	14	13	1	21	4	...	...	12	3500	5.0	15	...	...	...	...
51...	22	8	7	7	...	2	1	...	...	5	...	7.6	1	...	...	...	...
18-B	8802	2781	2358	2287	71	338	248	57	33	85	1823	5000	6.5	727	40	4.7	17
1	111	53	38	33	5	14	12	2	...	18	5000	6.6	25	...	...	...	...

# Isolate Table Body, Straighten Image

Rotate image around center of table – image is straightened

20

## City Block Characteristics

Table 2.—CHARACTERISTICS OF HOUSING UNITS, BY BLOCKS: 1960—Con.

["Total population" contains no persons in group quarters unless preceded by asterisk: one asterisk (\*) denotes less than 10 percent; two asterisks (\*\*), 10 percent or more]

Blocks within census tracts	Total population	All housing units by condition and plumbing							Occupied housing units								
		Total	Sound		Deteriorating			Dilapidated	Owner occupied			Renter occupied			Occu- pied by non- white	1.01 or more per sons per room	
			Total	With all plumb- ing facil- ities	Lack- ing some or all facil- ities	Total	With all plumb- ing facil- ities		Lacking some or all facilities	With flush toilet	No flush toilet	Total	Average value (dollars)	Aver- age num- ber of rooms			Total
21...	30	9	8	1	1	1	1	1	1	1	1	1	1	1	1	1	1
22...	87	28	22	6	11	5	1	1	1	1	1	1	1	1	1	1	1
23...	124	44	34	10	10	4	1	1	1	1	1	1	1	1	1	1	1
24...	**120	35	26	9	9	0	1	1	1	1	1	1	1	1	1	1	1
25...	247	55	38	17	14	3	1	1	1	1	1	1	1	1	1	1	1
26...	145	40	25	15	10	5	1	1	1	1	1	1	1	1	1	1	1
27...	85	20	12	8	7	1	1	1	1	1	1	1	1	1	1	1	1
28...	21	7	2	2	5	1	1	1	1	1	1	1	1	1	1	1	1
29...	14	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
30...	51	13	12	1	...	...	...	...	...	...	...	...	...	...	...	...	...
31...	18	4	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
34...	10	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
35...	9	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36...	54	18	7	11	8	1	1	1	1	1	1	1	1	1	1	1	1
37...	60	15	11	4	3	1	1	1	1	1	1	1	1	1	1	1	1
38...	132	52	38	14	8	6	1	1	1	1	1	1	1	1	1	1	1
43...	35	12	5	7	2	1	1	1	1	1	1	1	1	1	1	1	1
44...	6	2	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
45...	3	1	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
49...	141	43	17	16	26	20	3	3	1	1	1	1	1	1	1	1	1
50...	84	35	14	13	21	4	10	7	...	...	...	...	...	...	...	...	...
51...	22	6	7	7	2	1	1	1	...	...	...	...	...	...	...	...	...
18-B....	*8802	2781	2358	2287	71	338	248	57	33	85	1823	5000	6.5	727	40	4.7	171
1....	111	53	38	38	14	14	12	2	...	...	...	...	...	...	...	...	...

# Always Be Checking

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

21

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

22

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

23

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

24

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

21

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

22

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

23

Table 1 - CHARACTERISTICS OF FIBERING UNITS BY STACK (100-104)

Stack	Characteristics
100	...
101	...
102	...
103	...
104	...

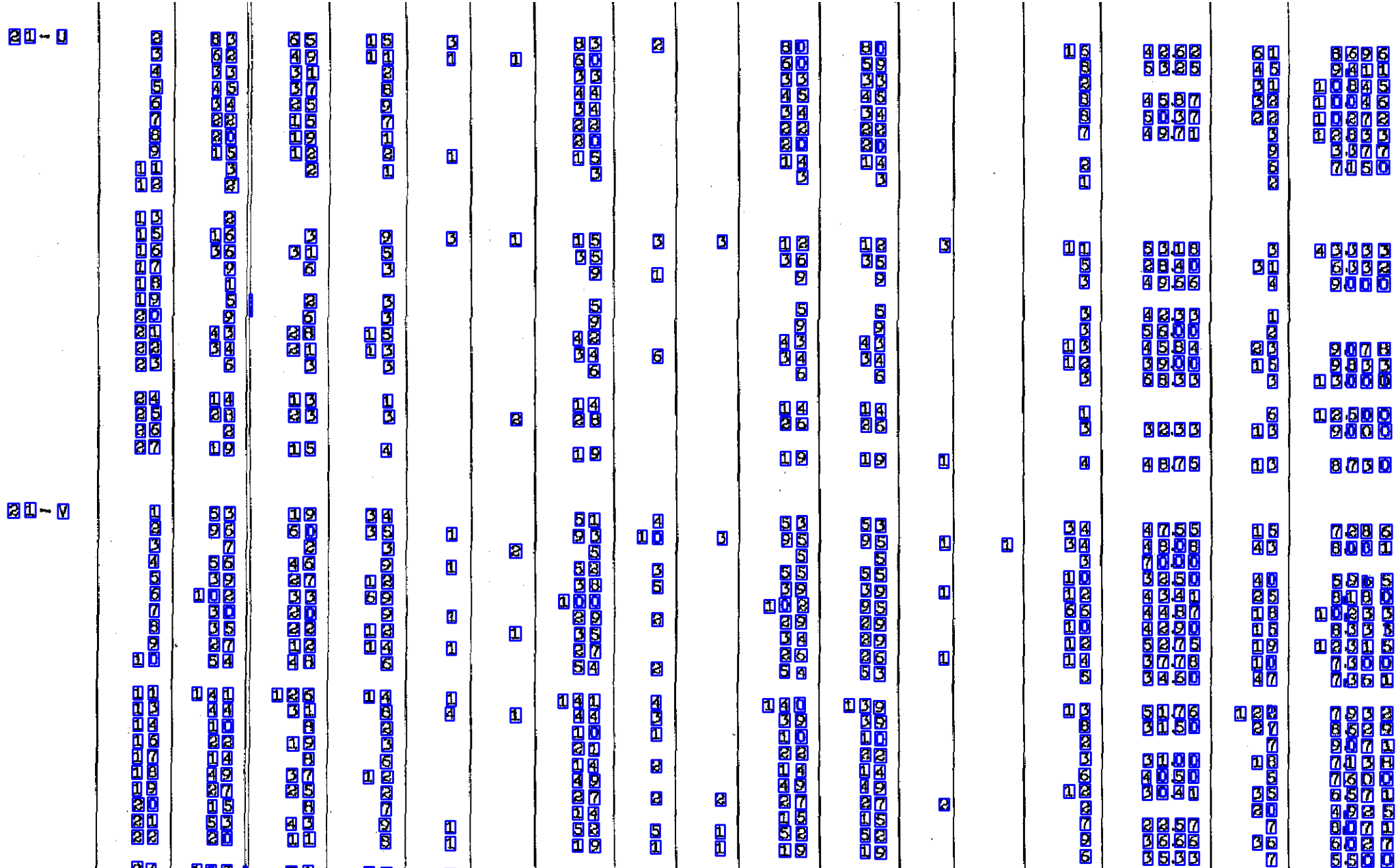
24





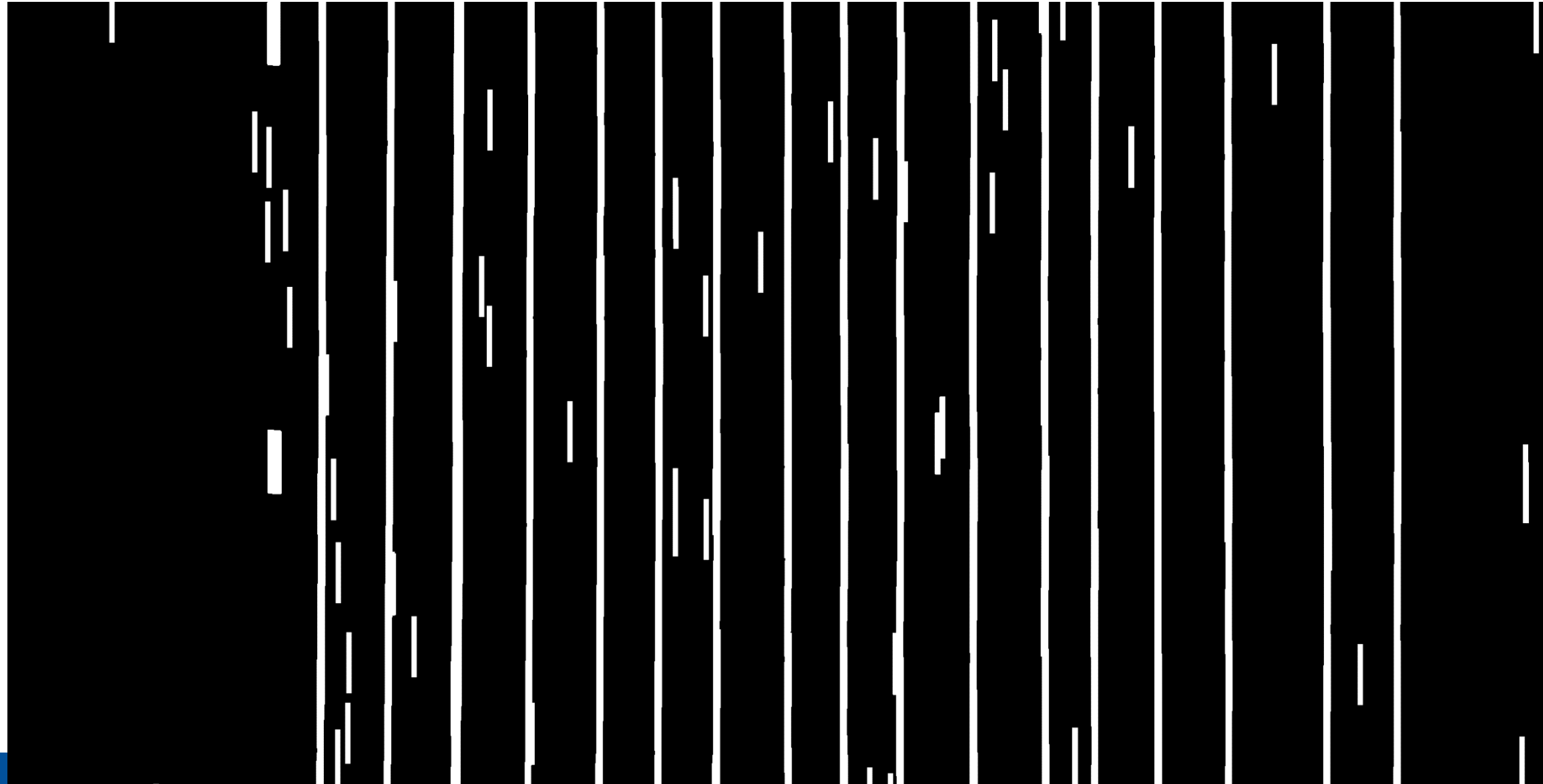
# Isolate Columns

Find everything that *could* be a character. Be aggressive, recall is important



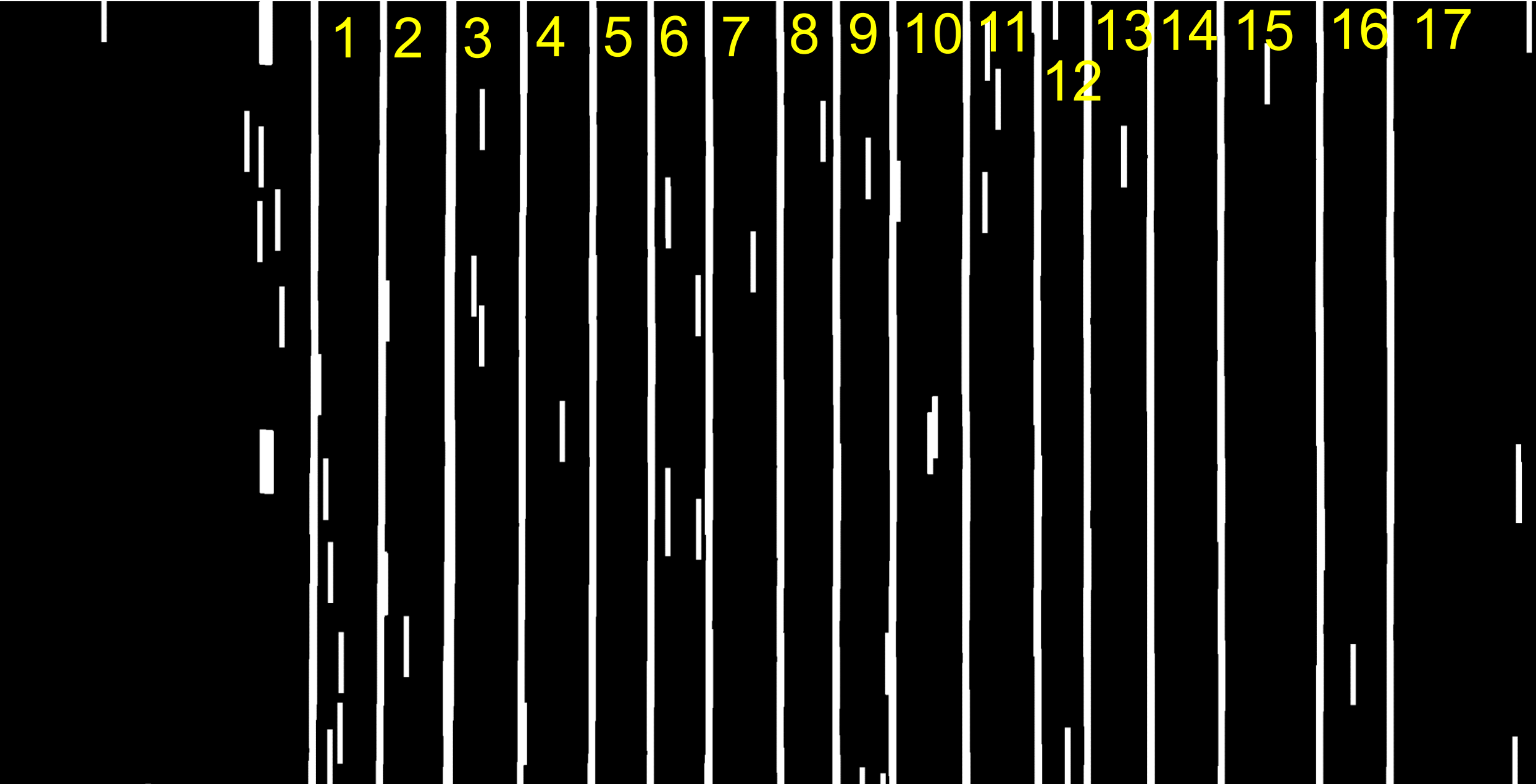
# Isolate Columns

Isolate and smear (slightly horizontally, aggressively vertically) what is left



# Isolate Columns

Find  $(N - 1)$  longest lines that are nearly vertical,  $N = \#$  of columns



# Isolate Columns

Columns are isolated

Row	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8	Col 9	Col 10	Col 11	Col 12	Col 13	Col 14	Col 15	Col 16	Col 17	Col 18	Col 19	Col 20
21-U	11	83	65	15	3	1	8	2	80	80	16	42.62	61	8.696						
	12	62	49	11	1		6		60	88	88	53.25	45	9.411						
	13	33	31	22			3		33	33	33		31	10.845						
	14	45	37	88			4		44	44	44	45.87	38	10.046						
	15	28	25	99			5		28	28	28	50.37	22	10.272						
	16	20	19	11			6		20	20	7	49.71	3	12.833						
	17	15	12	28			7		15	20	9		6	3.377						
	18	23	22	1			8		23	23	8		2	7.150						
	19	16	3	9			9		16	16	11	53.18	3	4.3333						
	20	36	31	35			3	3	36	36	5	28.40	31	6.332						
	21	9	6	3			1		9	9	3	49.66	4	9.000						
	22	5	2	3			2		5	5	3	42.33	1							
	23	4	2	1			3		4	4	3	56.00	2							
	24	3	1	3			4		3	3	1	45.84	23	9.078						
	25	6	3	3			5		6	6	1	39.00	15	9.833						
	26	14	13	1			6		14	14	3	68.33	3	13.000						
	27	28	23	3			7		28	28	1		6	12.500						
	28	2	2	3			8		2	2	3	32.33	13	9.000						
	29	19	15	4			9		19	19	4	48.75	13	8.730						
21-V	1	53	19	34			10		1	53	53	47.55	15	7.286						
	2	96	60	35			11		2	96	96	48.08	43	8.081						
	3	7	2	3			12		3	7	1	70.00								
	4	56	46	1			13		4	56	56	32.50	40	5.965						
	5	39	27	2			14		5	39	55	43.41	25	8.180						
	6	102	33	6			15		6	102	55	44.87	18	10.233						
	7	30	20	9			16		7	30	99	42.90	15	8.333						
	8	35	12	1			17		8	35	29	52.75	19	12.315						
	9	27	8	12			18		9	27	29	37.78	10	7.300						
	10	54	48	6			19		10	54	53	34.60	47	7.361						
	11	141	126	14			20		11	141	139	51.76	122	7.932						
	12	44	31	3			21		12	44	39	31.50	27	8.629						
	13	10	2	2			22		13	10	10		7	9.071						
	14	22	19	3			23		14	22	22	31.00	18	7.134						
	15	14	3	6			24		15	14	14	40.50	5	7.600						
	16	49	37	1			25		16	49	49	30.41	35	6.571						
	17	22	15	2			26		17	22	27	22.57	7	8.071						
	18	15	8	7			27		18	15	22	36.66	36	6.027						
	19	27	25	2			28		19	27	27		7							
	20	20	11	8			29		20	20	27		36	6.027						
	21	53	43	9			30		21	53	53	35.33	7	5.500						
	22	20	11	9			31		22	20	53		7							

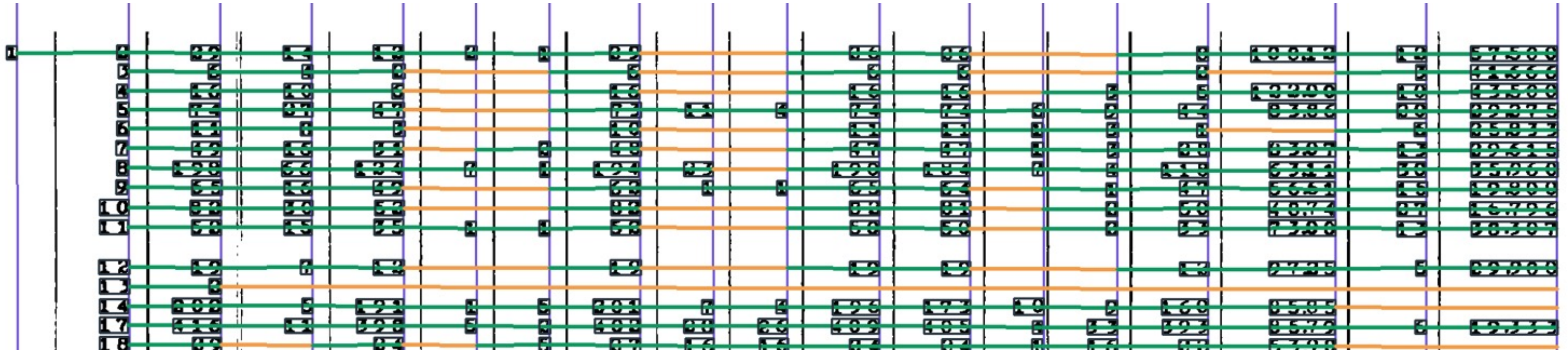
# Tesseract Each Row in Each Column

- Tesseract *highly* sensitive to input parameters, but flexible and governable
- Use restricted character set and character level confidence
- Collect character level text, bounding boxes and confidence

# Use Table's Internal Structure to Build Rows and Columns

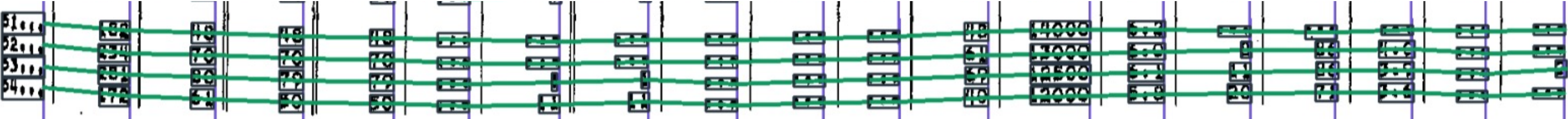
- Start with block column (always populated)
- Look right to find the two-way unique nearest neighbor for each row, requiring the angle to the nearest neighbor be similar for all rows
- Create a synthetic cell for cells that do not have a nearest neighbor conforming to angle and distances of other cells in column
- Repeat moving out left and right to cover all columns
- Create PDF of all pages to scan for errors

# Use Table's Internal Structure to Build Rows and Columns





# Use Table's Internal Structure to Build Rows and Columns



# Use Table's Internal Structure to Build Rows and Columns



# Train and Apply Custom Model

- Match cells to training data – Washington DC, Mapping Segregation
- Train random forest model at the character level
  - Pixel value by position in bounding box
  - Tesseract predicted text
  - Tesseract confidence
- Grid search with cross validation to tune hyperparameters
- Apply model (out of sample) to remaining cities

# Identify Internal Inconsistencies, Compare to Tract Totals

- Internal consistency, e.g. Owner Occupied + Renter Occupied = Occupied
- Check for outliers at column level
- Compare stats to tract totals, accounting for suppression
- Make corrections easy with Excel tool

tract	ocr	human
<b>T O T A L</b>	TOTAL	
<b>1 - A</b>	1-A	
<b>1 - B</b>	1-B	
<b>1 - C</b>	1-C	
<b>2 - A</b>	2-A	
<b>2 - B</b>	2-B	
<b>2 - C</b>	2-C	
<b>3 - A</b>	3-A	
<b>3 - B</b>		

# Caveats

- Approach requires some customization per dataset
- Manual steps remain (and probably always will)
  - Identifying unusable scans
  - Identification of page ranges in source documents (missing pages)
  - Always be checking
  - Tract transcription is still manual

# Current State

- Scaling work to all 16 cities for 1950
- Refining issues with 1960 model
- Starting 1940 work
- Textract for assist with tract identifiers?
- Claude or other LLM based service for first cut?

# Summary

# Summary

- We are working on **digitizing** the historical Censuses of Housing **Block Statistics**, 1940 to 1970.
- Our goal: Develop & release data for 16 cities, training & validation data, and methods & code.
- The three major tasks are digitizing block **shapes**, the block **situations**, and the block **statistics**.
- This is a work in progress; Questions and comments welcome!



Thanks!