

Generating Useful Data with Computer Vision Tools: 2 Use Cases (PDFs & Imagery)



Location matters

Erik Neemann
10 May 2023

Overview



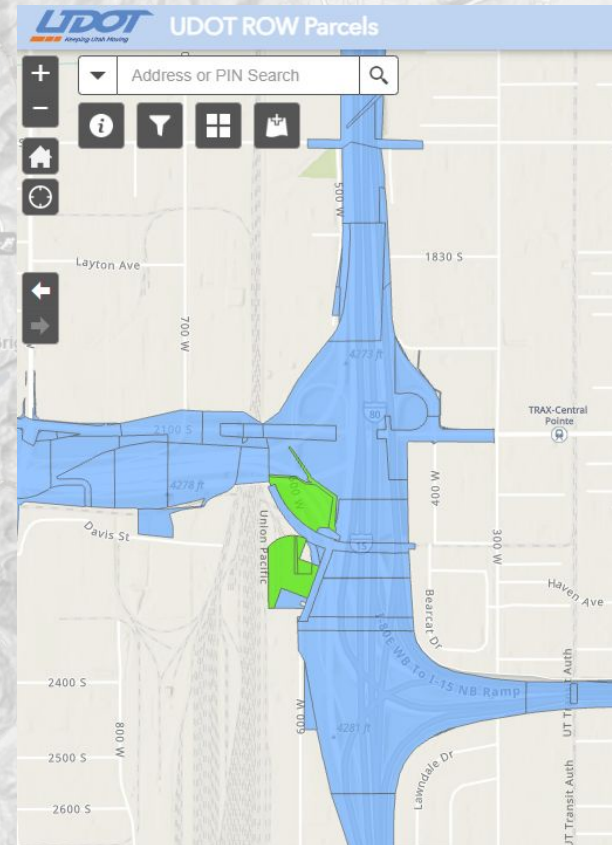
- **Two Computer Vision (CV) Projects**
 - **UDOT Parcel Detection Project**
 - **DHHS Cooling Towers Project**
- **Motivation**
- **Process**
- **Details/Tools**
- **Results**



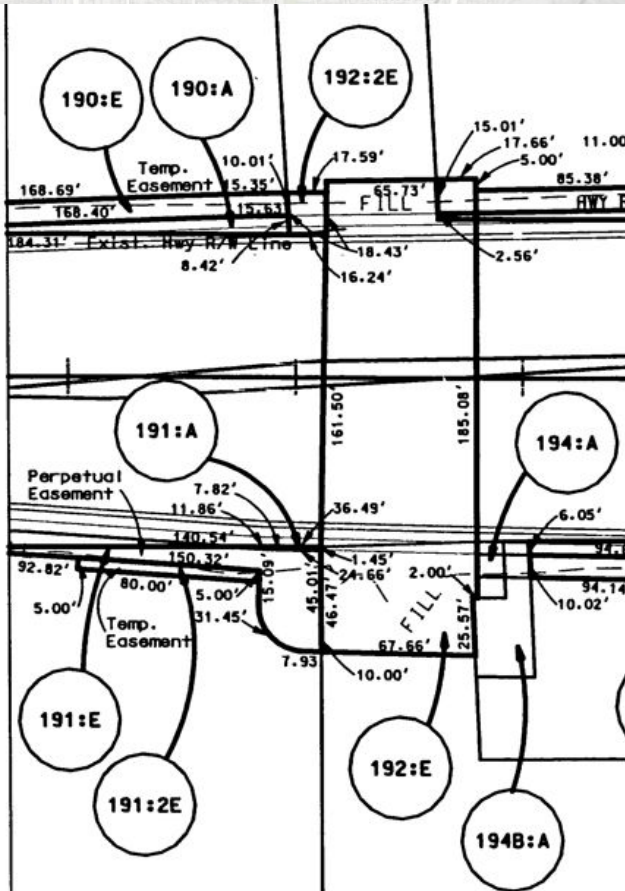
Utah Department of
Health & Human
Services

Motivation - UDOT Parcel Detection

- UDOT has acquired a LOT of property over 100+ years
 - ...but they didn't keep track of where that property was located
- UDOT has a tons (90K) of project plan documents
 - ...but it's difficult to find a specific parcel within those documents
- How can they untangle where everything is?
 - Manually sift through 90,000 documents?
 - NO!
 - Let the machines do the work!



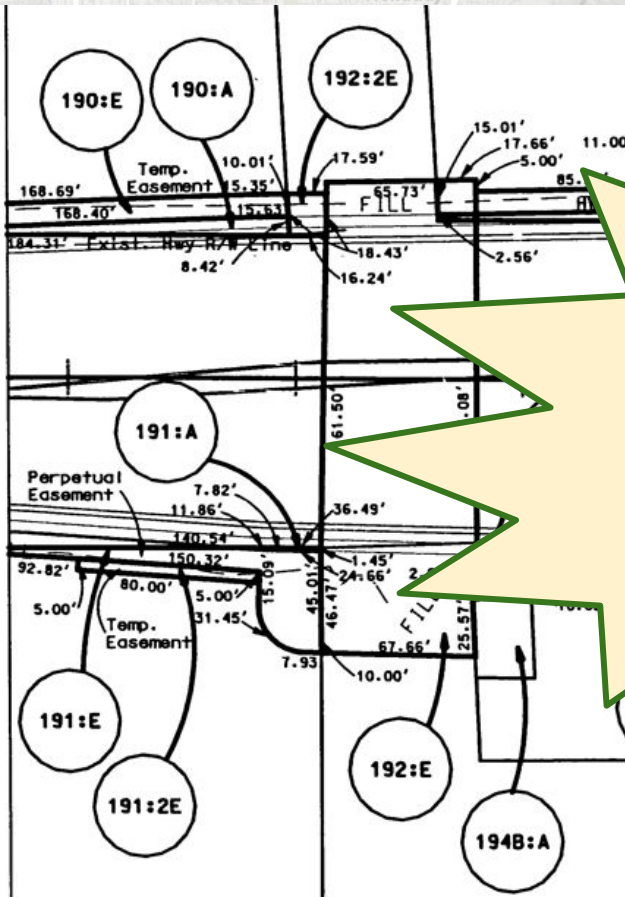
Motivation - UDOT Parcel Detection



- Parcel information is annotated in a very specific way - circles around the text
- Plan documents fairly consistent in format
- Circles take up roughly the same size in any given document
- Parcel text follows patterns, w/ defined rules
 - numbers/letters, colon, numbers/letters
 - 193B:2A
 - 191:E
 - 41BNT:2E



Motivation - UDOT Parcel Detection



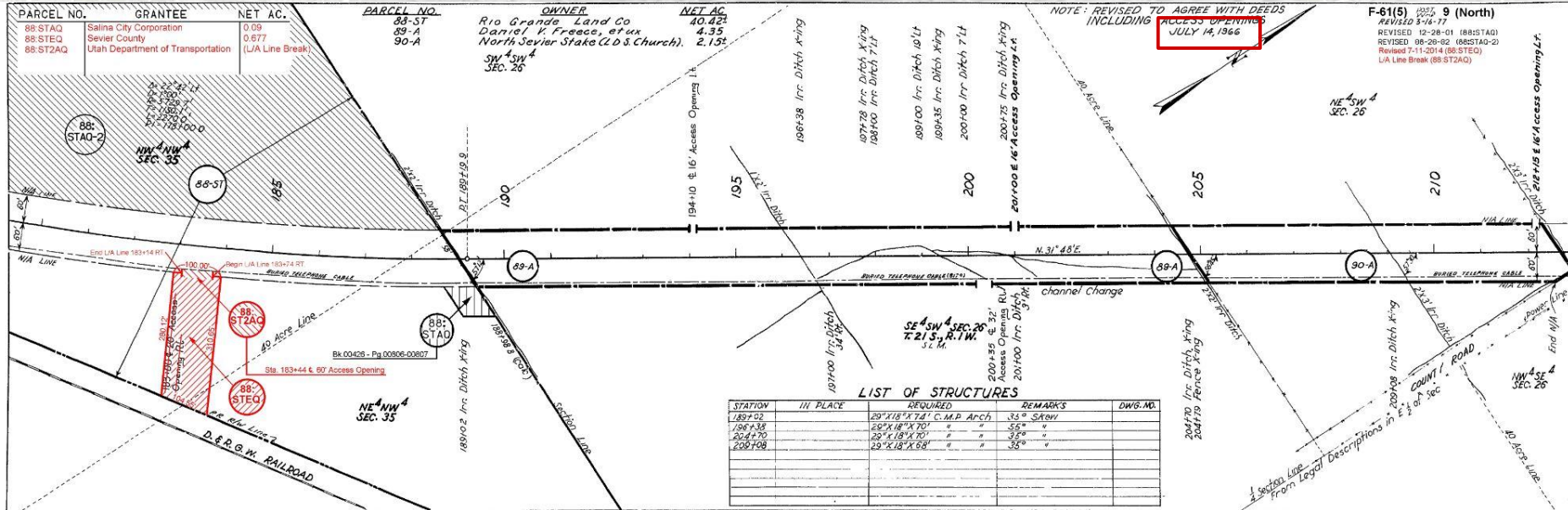
- Parcel information is annotated in a very specific way - circles around the text
- Plan annotations are fairly consistent in format and the same size in any

This problem is solvable with computer vision!

- 191:E
- 41BNT:2E



Examples



COUNTY OF KANE, ARIZONA
 COURTESY OF THE ARIZONA DEPARTMENT OF LAND AND WATER
 BY: [Name] PLS. 133
 DATE: [Date]
 SHEET NO. 9
 TOTAL SHEETS: 9

SARATOGA SPRINGS CITY
(INCORPORATED)



PARCEL NO.	OWNER	NET AC.	SQ. FEET	EXIST./R/W AC. IN DEED	REMAINING AC. LEFT	REMAINING AC. RIGHT
0182: 808:A	COLLINS BROTHERS LAND DEVELOPMENT LLC	1.168	50.860	NONE	143.258	NONE
0182: 808:2A	COLLINS BROTHERS LAND DEVELOPMENT LLC	7.259	316.219	0.540	135.999	NONE

5/11/2018 1:36:01PM

Z:\000\00072 M/C Redwood Road St. County\Information\Maps\VT03_MF-08066_RW-34_Aerial.dwg

REVISIONS

NO.	DATE	APPROVED BY	JPA	REMARKS

DRAWN BY: DBA
 QC CHECKED BY: JPA
 DATE: _____

RIGHT OF WAY ENGINEER: _____

UTAH DEPARTMENT OF TRANSPORTATION
 REGION 2 - SALT LAKE CITY, UTAH

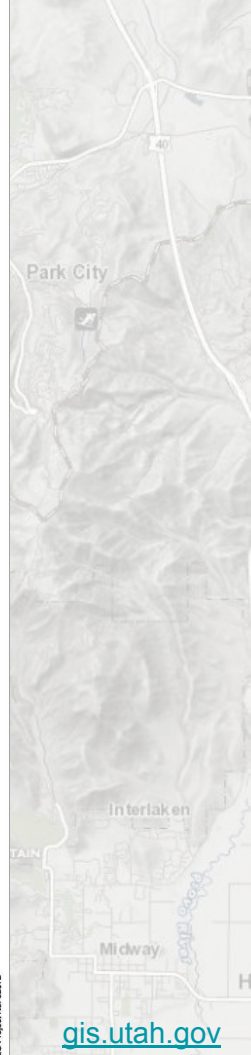
APPROVED: _____

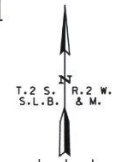
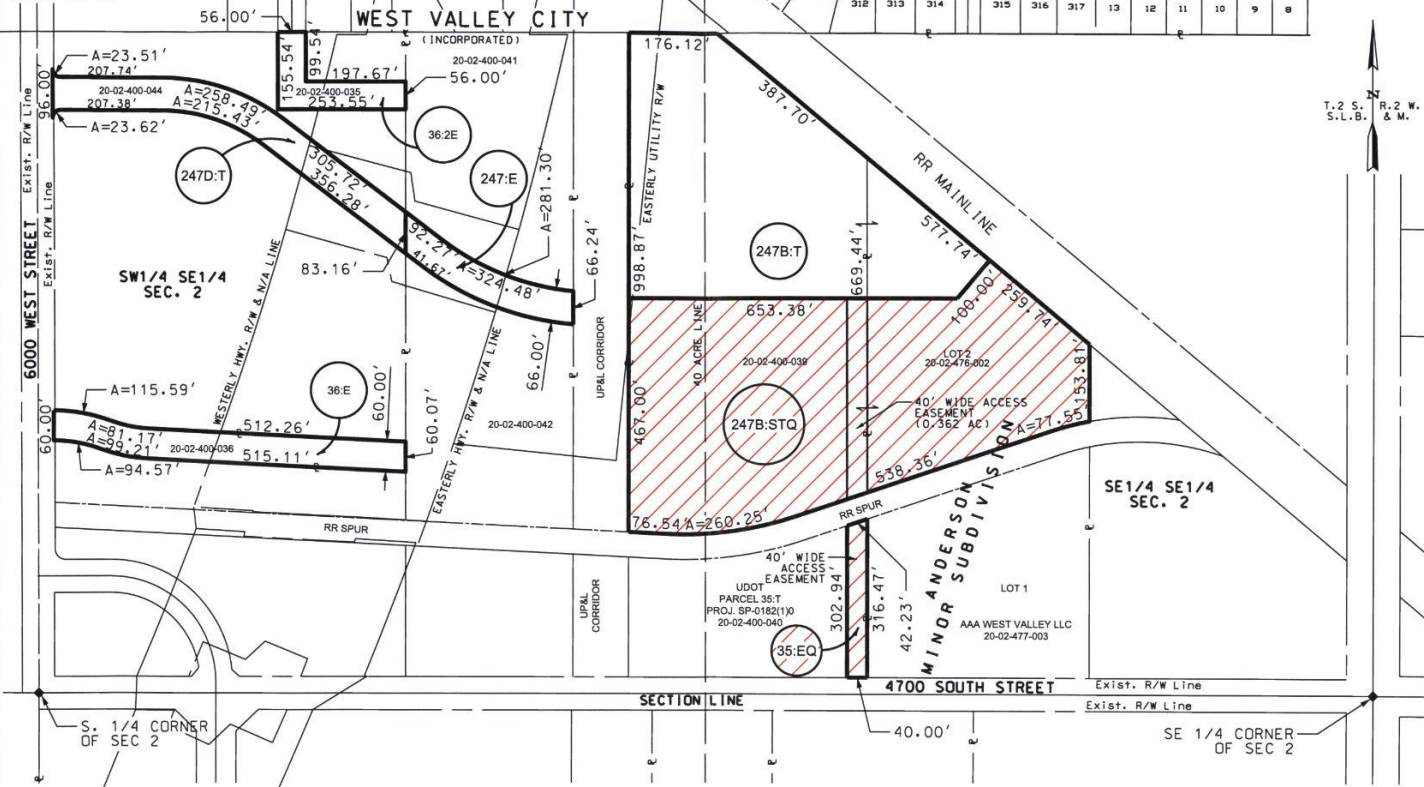
MOUNTAIN VIEW CORRIDOR
 REDWOOD RD - 9000 SOUTH
 RIGHT OF WAY PLAN SHEET

PROJECT NUMBER: SP-0182(1)0

TEMPORARY SHEET NO: RW-34

8000 Project No. 00072





PARCEL NO.	OWNER	NET AC.	SQ. FEET	EXIST. R/W AC. IN DEED	OWNERSHIP AC.	REMAINING AC. LEFT	REMAINING AC. RIGHT
0182:247B:T	BB SOLD PC	13.962	608,185	NONE	13.962	NONE	NONE
0182:247D:T	BB SOLD PC	1.200	52,274	NONE	1.200	NONE	NONE
0182:247E	BB SOLD PC	0.560	24,411	NONE	NONE	RELEASE OF EASEMENT	NONE
0182:36:E	BB SOLD PC	0.977	42,537	NONE	NONE	RELEASE OF EASEMENT	NONE
0182:36:2E	GIOVENGA PROPERTIES, LLC	0.454	19,776	NONE	NONE	RELEASE OF EASEMENT	NONE

PARCEL NO.	GRANTEE	NET AC.
0182:247B:STQ		8.091
0182:35:E		0.284



UTAH DEPARTMENT OF TRANSPORTATION
RIGHT OF WAY DESIGN

PROJECT: MOUNTAIN VIEW CORRIDOR
PRODUCT NUMBER: SP-0182(1)0
SHEET NO: TEMP 54

DATE: 3/08/11

APPROVED: [Signature]
PROFESSIONAL LAND SURVEYOR

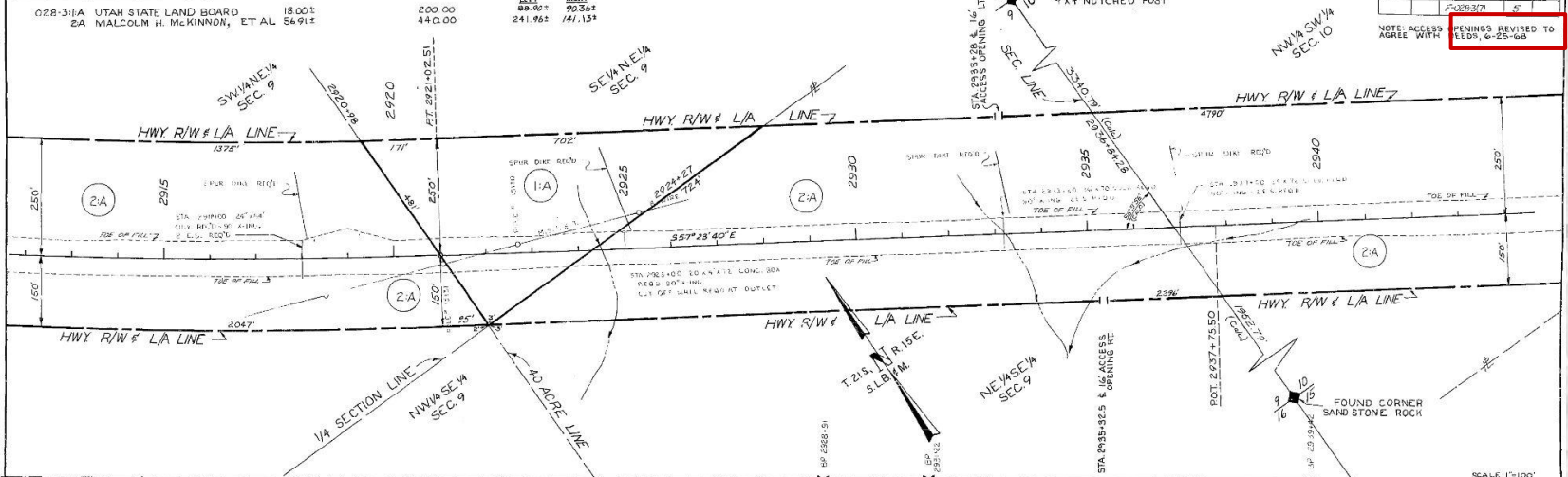
DRAWN BY: CRH
CHECKED BY: JCM
DATE: 3/08/11



PARCEL NO	OWNER	NET AC.	SQ. FT.	OWNERSHIP AD.	REMAINING AC.
028-31A	UTAH STATE LAND BOARD	18.00±	200 00		
2A	MALCOLM H. MCKINNON, ET AL	56.91±	44 00 00		
				LEFT	RIGHT
				888.96±	90.56±
				241.96±	141.15±

US 89 TO GREEN RIVER TOWARDS WOODSIDE

		4286370	57
NOTE: ACCESS OPENINGS REVISIONS REVISED TO 6/25/16			

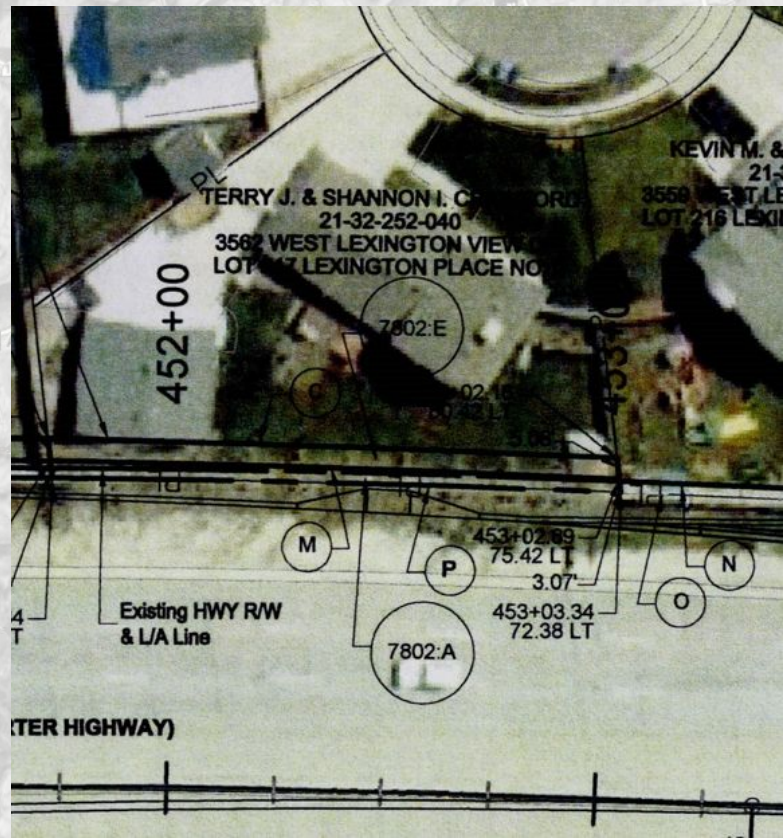


ELEVATION	EMPHASIS		SECTION		SECTION		SECTION		SECTION		SECTION		SECTION		SECTION		SECTION		SECTION		SECTION	
	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	
8250																						
70																						
60																						
5250																						



UDOT Parcels - Challenges

- UDOT completed a bulk export of files from their project software into Google Cloud Storage
 - ...but many files are not relevant!
- Various file types (documents, images, other files)
- Images vary in format and file type
- Images vary in size, resolution, DPI
- Images vary in font style (typed, handwritten, different typed fonts)
- Images vary in orientation
- Images vary in quality, noise, consistency



Examples



MILESTONE DESIGN TRANSMITTAL

Legacy Parkway Design Build Project

PLAN PACKAGE (if applicable)

Segment: **1-15** Type: Sequence:

Original DocC Number: **23DCN87**

Title: **Design Change Notice**

SUBJECTS (Check box for each category that defines transmittal contents) X

Roadway	Drainage	Geotechnical	Pavement	Structures
MOT	Signals	Lighting	Striping	Signs
ATMS	Utility	Landscape	Aesthetics	
Other:				

REVIEW ITEMS
Use additional sheet to list items

1	Revision 1 to Phase 1 Lighting Plans	
2	Check prints	
3	DCN Form	
4		<i>NO COMMENTS</i>
5		
6		

IQF FOR AUDIT

1st DocC Extension #: **23DCN87-1**

Transmitted By: Dr. K. N. Gunalan
FAK-LLC Design Manager
360 N. 700 W. Suite F
North Salt Lake, UT 84054

Transmitted To: Mr. John Bale
IQF Manager
360 N. 700 W. Suite F
North Salt Lake, UT 84054

K. N. Gunalan
Signature 10/16/03
Date

John Bale
Signature 10/22/03
Date

Ben Audit
Print IDQM Signature Date

Ben Audit
Print IDQM Signature Date

IDQM sign and return (with comments) to FAK Design Manager

UDOT FOR REVIEW

2nd DocC Extension #: **23DCN87-2**

Transmitted By: Dr. K. N. Gunalan
FAK-LLC Design Manager
360 N. 700 W. Suite F
North Salt Lake, UT 84054

Transmitted To: Mr. Todd Jensen
UDOT Design Manager
360 N. 700 W. Suite F
North Salt Lake, UT 84054

K. N. Gunalan
Signature 10-22-03
Date

Todd Jensen
Signature 11/4/03
Date

Beverly Smith
Print UDOT Oversight Signature Date

UDOT sign and return (with comments) to FAK Design Manager within five (5) days of receipt. UDOT signature demonstrates concurrence with the subject of this transmittal when applied to the specifically stated subject. Signature is not a Release for Construction.
FAAdmin\Forms\Milestone Transmittal.doc

Fluor Ames Kraemer, LLC

TRANSMITTAL

No. 01892

360 North 700 West, Ste F
North Salt Lake City, UT 84054

Phone: 801-951-1900
Fax: 801-951-1840

PROJECT: Legacy Parkway Design-Build

DATE: 5/13/2004

TO: Utah Department of Transportation
360 North 700 West, Ste F
North Salt Lake, UT 84054

REF: Plan Set, Park and Main,
final revision, from CRS
AFC

ATTN: *Rick Campagna*

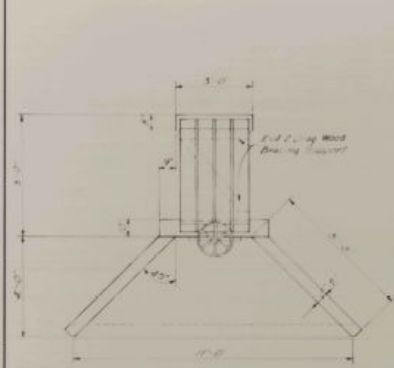
WE ARE SENDING:	SUBMITTED FOR:	ACTION TAKEN:
<input type="checkbox"/> Shop Drawings	<input type="checkbox"/> Approval	<input type="checkbox"/> Approved as Submitted
<input type="checkbox"/> Letter	<input checked="" type="checkbox"/> Your Use	<input type="checkbox"/> Approved as Noted
<input type="checkbox"/> Prints	<input checked="" type="checkbox"/> As Requested	<input type="checkbox"/> Returned After Loan
<input type="checkbox"/> Change Order	<input type="checkbox"/> Review and Comment	<input type="checkbox"/> Resubmit
<input type="checkbox"/> Plans		<input type="checkbox"/> Submit
<input type="checkbox"/> Samples	SENT VIA:	<input type="checkbox"/> Returned
<input type="checkbox"/> Specifications	<input checked="" type="checkbox"/> Attached	<input type="checkbox"/> Returned for Corrections
<input type="checkbox"/> Other	<input type="checkbox"/> Separate Cover Via: hand	<input type="checkbox"/> Due Date:

ITEM	NUMBER	REV. NO.	DESCRIPTION	STATUS
001	6	5/13/2004	Plan Set Park & Main, Final Revision from CRS	AAP

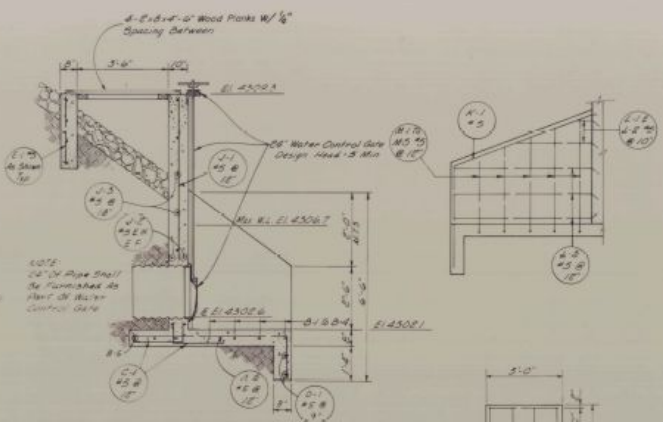
Remarks:
*58M.
1 LG.*

CC: file 220, John Bale, Alan Beane, Doc Cntrl North, Cliff Barber, Lewis Young, Jerome Frank, Bryan Keck, Dan Openshaw

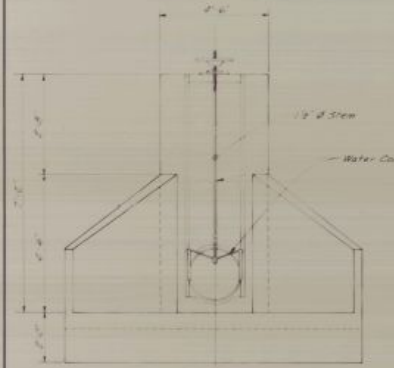
Signed: *Deena Farmer*
Deena Farmer



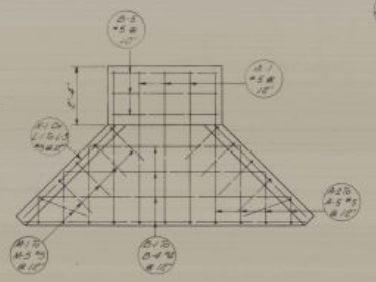
PLAN



SECTION



ELEVATION



SLAB PLAN

INLET STRUCTURE, STATION 765+75 RT.

LOCATION	MATERIAL	SIZE	QTY	FACTOR	TOTAL QUANTITY	REMARKS
SLAB	A-1	S	3	7'-8"	23'-6"	a = 5'-10"
	A-2	S	2	7'-8"	15'-4"	a = 5'-10"
	A-3	S	2	4'-8"	8'-0"	a = 3'-0"
	A-4	S	2	3'-8"	7'-0"	a = 2'-0"
CLB	B-1	S	1	11'-2"	11'-2"	a = 1'-0"
	B-2	S	1	3'-8"	3'-8"	
	B-3	K	1	7'-8"	7'-8"	
	B-4	S	1	3'-8"	3'-8"	
	B-5	S	3	4'-8"	14'-0"	
SLAB	C-1	S	8	2'-6"	20'-0"	
CUT OFF WALL	D-1	S	2	11'-2"	22'-4"	
RAILING	E-1	S	8	2'-8"	21'-6"	
WIND WALL	F-1	S	2	9'-4"	18'-8"	
WIND WALL	F-2	S	2	6'-0"	12'-0"	
WIND WALL	F-3	S	2	4'-2"	8'-4"	
WIND WALL	G-1	S	2	0'-8"	1'-6"	
WIND WALL	H-1	S	2	3'-0"	6'-0"	
WIND WALL	H-2	S	2	5'-8"	11'-6"	
WIND WALL	H-3	S	2	6'-8"	13'-6"	
WIND WALL	M-1	S	2	3'-11"	7'-2"	
WIND WALL	M-2	S	2	5'-7"	11'-4"	
WIND WALL	M-3	S	2	6'-0"	12'-0"	
WIND WALL	M-4	S	2	6'-4"	12'-8"	
TOTAL					395'-8" @ 1043 $\frac{1}{4}$ = 413'	

ESTIMATED QUANTITIES

- Concrete "Class 24" (AF) **3,533** Cu Yds
- Reinforcing Steel **413** Lbs
- 24" Water Control Gate, Type I **2000**
- Timber Untreated **0030** 1000 BB FT

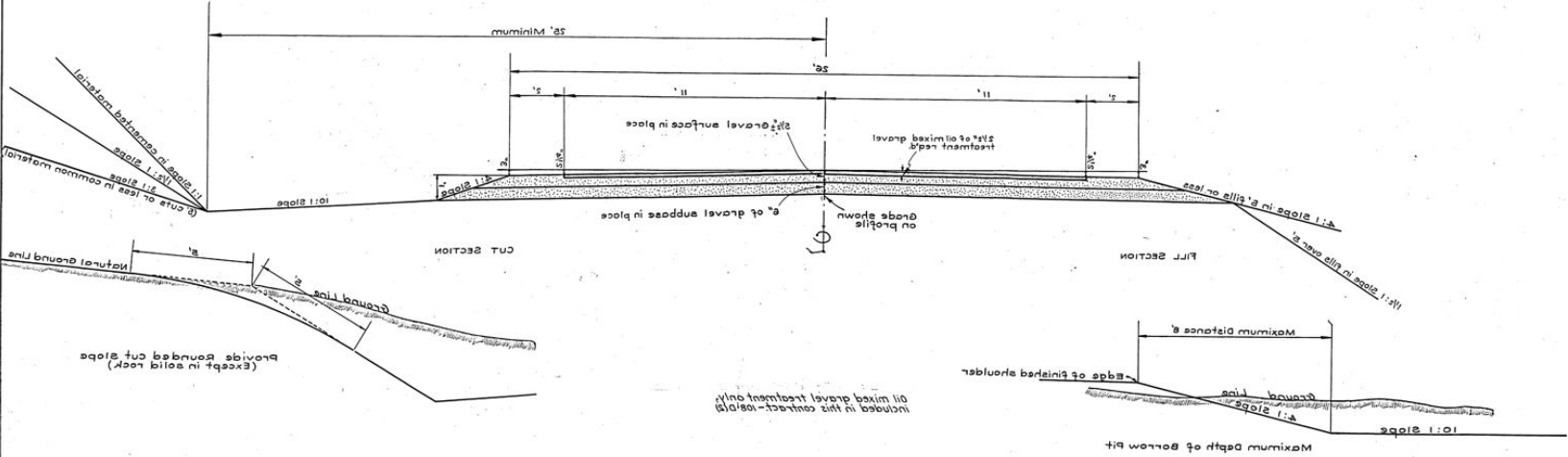
NOTES
 All Structural Steel Shall Be Painted In Accordance With The Utah Department Of Highway Standard Specifications, Dates 1970 And Revisions Thereof.

UTAH DEPARTMENT OF TRANSPORTATION			
SALT LAKE CITY, UTAH			
ROADWAY DIVISION			
LAGOON TO LAYTON			
IRRIGATION UTILITY			
DESIGN C.H.S.	76	CHECK C.S.	JG
DESIGN P.E.	LS	CHECK C.H.S.	DL
DESIGN C.H.S.	LS	CHECK I.I.S.	LS
DATE	11-1-76	BY	J.D. Brown
APPROVED	10/1/76	BY	J.D. Brown
DAVIS BRADY			

NO.	BY	DATE	TYPE	REMARKS

TYPICAL CROSS SECTION

PROJECT NO.	108-0103	DATE	5-2-39
DISTRICT	UTAH	SCALE	AS SHOWN
SECTION	2	DESIGNED BY	W. J. WILSON
DATE	5-2-39	CHECKED BY	W. J. WILSON



Provide Rounded cut slope (Except in solid rock)

Oil mixed gravel treatment only included in this contract - (a)(b)

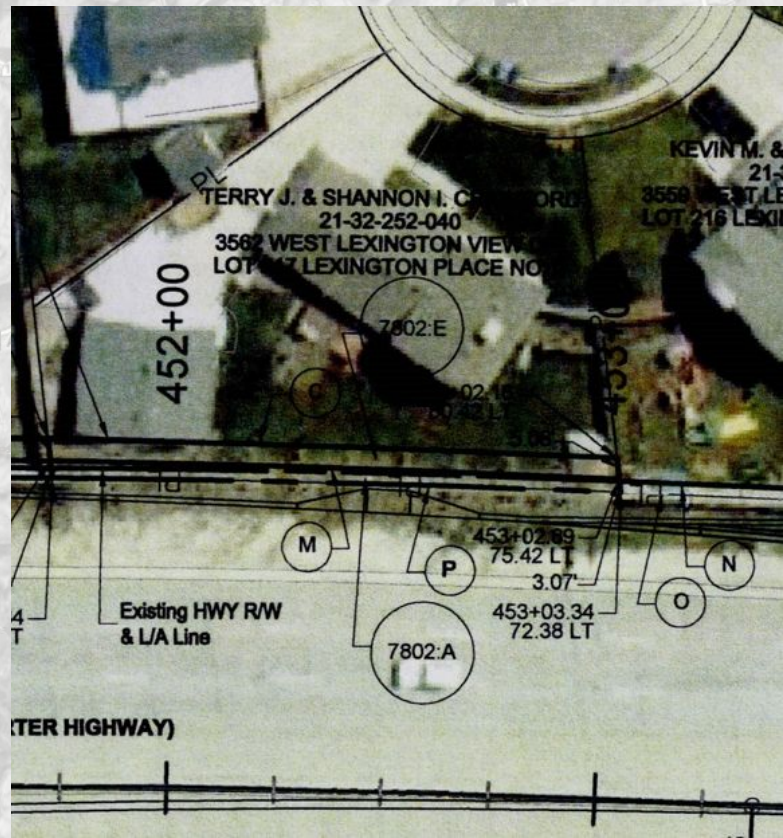
Work started on APR 108-0103 4-24-39
 Completed 5-2-39

UTAH STATE ROAD COMMISSION
 52' x 2 1/2" OIL MIXED GRAVEL SURFACE ROAD
 TYPE --
 Revised 4-17-39

Widen and super-elevate curves according to 1-24-39

UDOT Parcels - Challenges

- UDOT completed a bulk export of files from their project software into Google Cloud Storage
 - ...but many files are not relevant!
- Various file types (documents, images, other files)
- Images vary in format and file type
- Images vary in size, resolution, DPI
- Images vary in font style (typed, handwritten, different typed fonts)
- Images vary in orientation
- Images vary in quality, noise, consistency



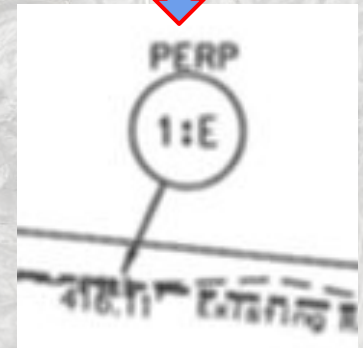
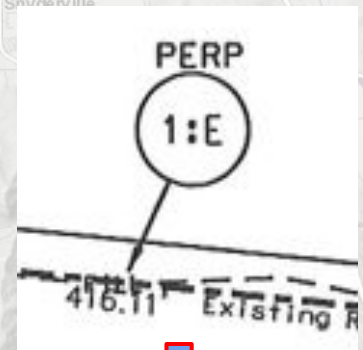
Step 1 - Initial File Processing

- About 90K objects in cloud storage bucket
- Several files were irrelevant and discarded (.xlsx, .doc, others)
- Many files types to deal with
 - ~46K PDF Documents (.pdf)
 - ~42K images (.tif, .jpg)
- Convert PDFs to images
 - Multipage PDFs to multiple images



Step 2 - Detect Circles

- Image preprocessing to improve circle detection
 - Convert to grayscale
 - Add slight blur
- Detect Circles
 - Assume circle radius ~2.5% of image width
 - Iterate through up to 6 values (smaller, bigger, smaller, etc.)
 - Stop when 1-100 circles are found
- Build Mosaics
 - Crop out square around detected circle
 - Mask out area outside circle
 - Inset the mask to remove circle outline
 - Stitch together all cropped squares into a mosaic
- Upload to Google Cloud Storage
- 10,000(!) worker tasks running in parallel in Google Cloud Run



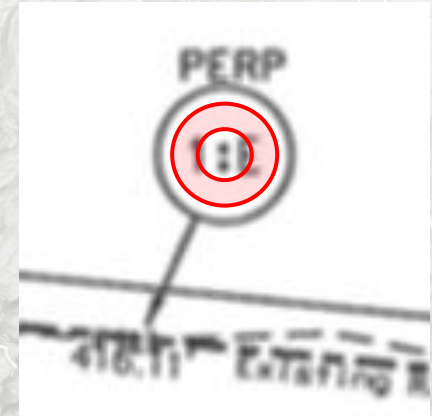
Grayscale & Blur

- 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
- 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**

Step 2 - Detect Circles

- Image preprocessing to improve circle detection
 - Convert to grayscale
 - Add slight blur
- Detect Circles
 - Assume circle radius $\sim 2.5\%$ of image width
 - Iterate through up to 6 values (smaller, bigger, smaller, etc.)
 - Stop when 1-100 circles are found
- Build Mosaics
 - Crop out square around detected circle
 - Mask out area outside circle
 - Inset the mask to remove circle outline
 - Stitch together all cropped squares into a mosaic
- Upload to Google Cloud Storage
- 10,000(!) worker tasks running in parallel in Google Cloud Run
 - 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
 - 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**

radius too small

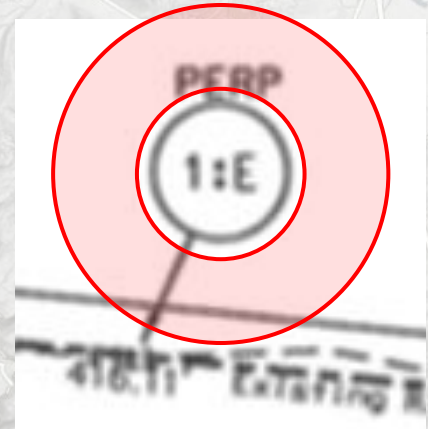


Iterate Through
Radius Sizes

Step 2 - Detect Circles

- Image preprocessing to improve circle detection
 - Convert to grayscale
 - Add slight blur
- Detect Circles
 - Assume circle radius $\sim 2.5\%$ of image width
 - Iterate through up to 6 values (smaller, bigger, smaller, etc.)
 - Stop when 1-100 circles are found
- Build Mosaics
 - Crop out square around detected circle
 - Mask out area outside circle
 - Inset the mask to remove circle outline
 - Stitch together all cropped squares into a mosaic
- Upload to Google Cloud Storage
- 10,000(!) worker tasks running in parallel in Google Cloud Run
 - 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
 - 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**

radius too big



Iterate Through
Radius Sizes



Step 2 - Detect Circles

- Image preprocessing to improve circle detection

- Convert to grayscale
- Add slight blur

- Detect Circles

- Assume circles are small
- Iterate over all pixels
- Stop when no more circles are found

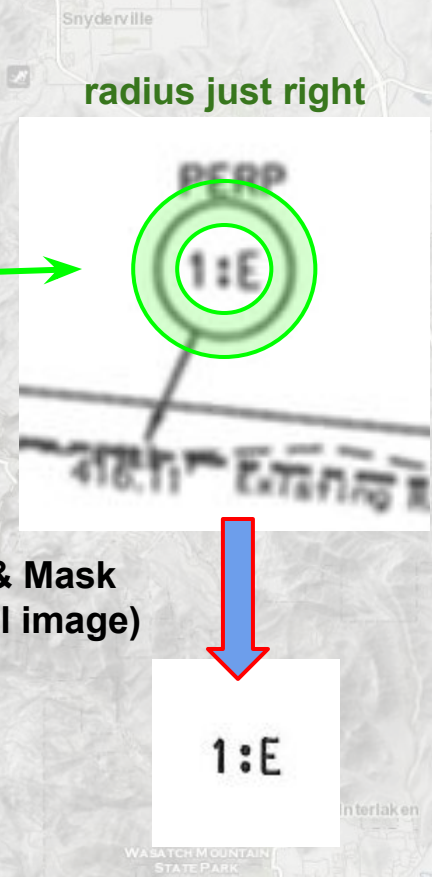
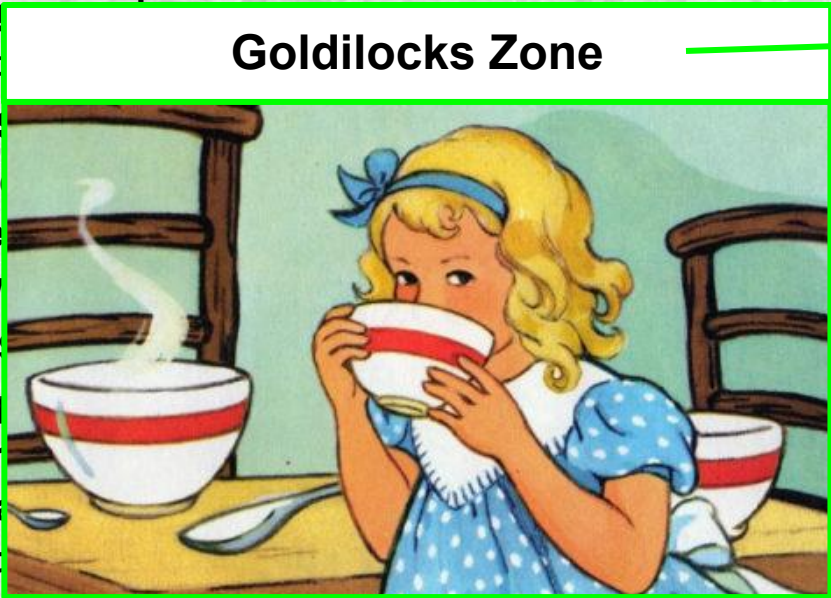
- Build Mosaic

- Crop out small circles
- Mask out areas with no circles
- Inset the mosaic
- Stitch together

- Upload to Google Cloud Storage

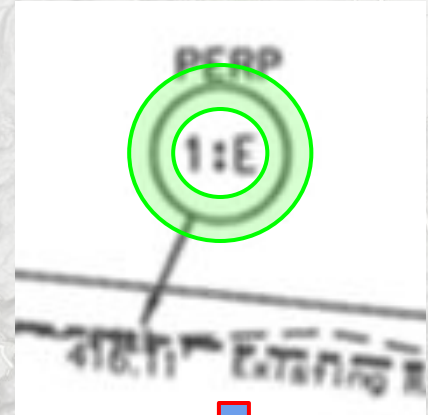
- 10,000(!) worker tasks running in parallel in Google Cloud Run

- 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
- 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**



Step 2 - Detect Circles

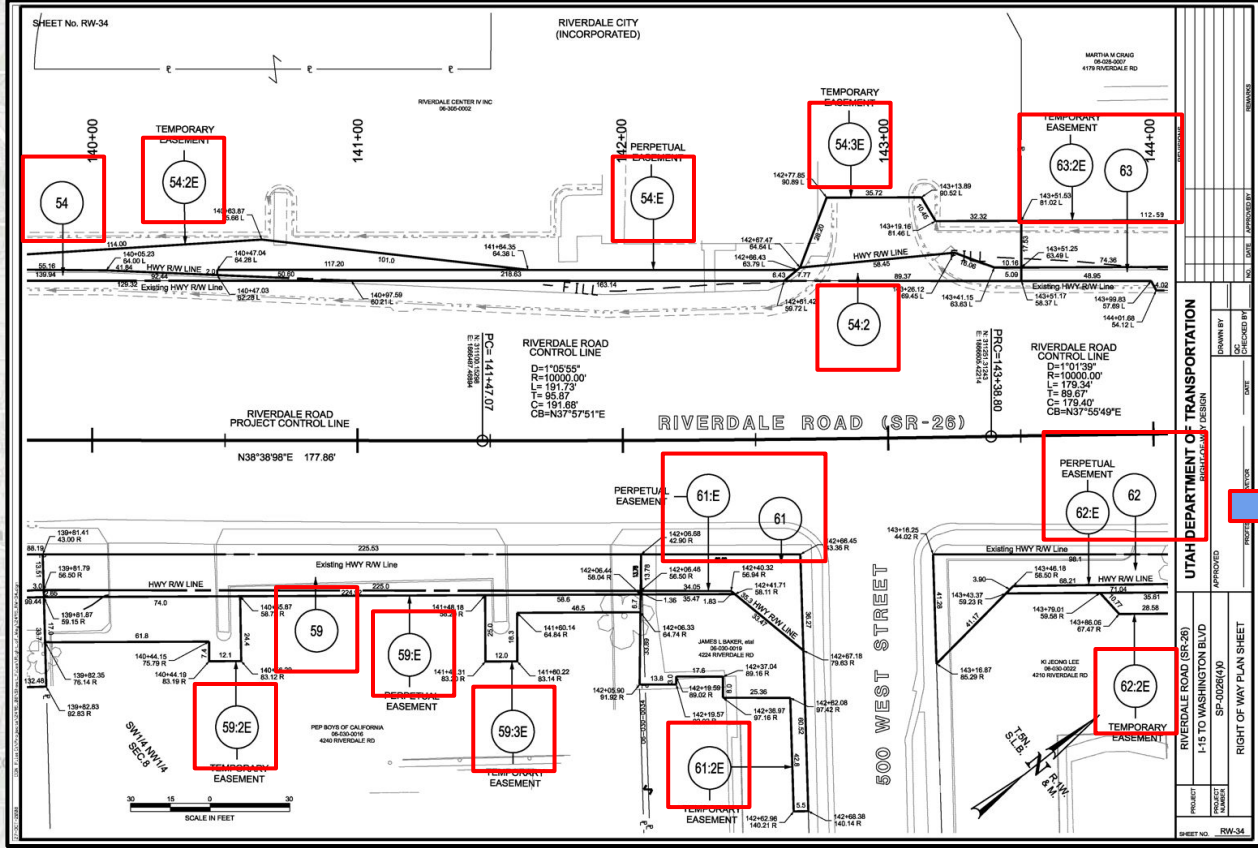
- Image preprocessing to improve circle detection
 - Convert to grayscale
 - Add slight blur
- Detect Circles
 - Assume circle radius ~2.5% of image width
 - Iterate through up to 6 values (smaller, bigger, smaller, etc.)
 - Stop when 1-100 circles are found
- Build Mosaics
 - Crop out square around detected circle
 - Mask out area outside circle
 - Inset the mask to remove circle outline
 - Stitch together all cropped squares into a mosaic
- Upload to Google Cloud Storage
- 10,000(!) worker tasks running in parallel in Google Cloud Run
 - 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
 - 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**



Crop & Mask
(original image)



Step 2 - Detect Circles

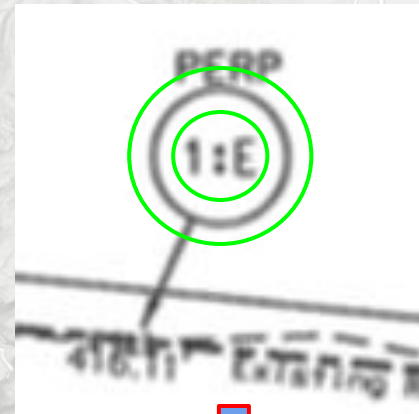


59:2E	62:2E	54	63:2E
54:E	63	59	62
61	54:2	59:E	59:3E
54:3E	61:E	54:2E	62:E
61:2E			

Build Mosaic

Step 2 - Detect Circles

- Image preprocessing to improve circle detection
 - Convert to grayscale
 - Add slight blur
- Detect Circles
 - Assume circle radius ~2.5% of image width
 - Iterate through up to 6 values (smaller, bigger, smaller, etc.)
 - Stop when 1-100 circles are found
- Build Mosaics
 - Crop out square around detected circle
 - Mask out area outside circle
 - Inset the mask to remove circle outline
 - Stitch together all cropped squares into a mosaic
- Upload to Google Cloud Storage
- 10,000(!) worker tasks running in parallel in Google Cloud Run
 - 2.5 run of 10K or 25K, each took about 3 minutes = 75,000 minutes of processing
 - 52 days or **7.4 weeks of processing time completed in less than 4 hours!!!!**

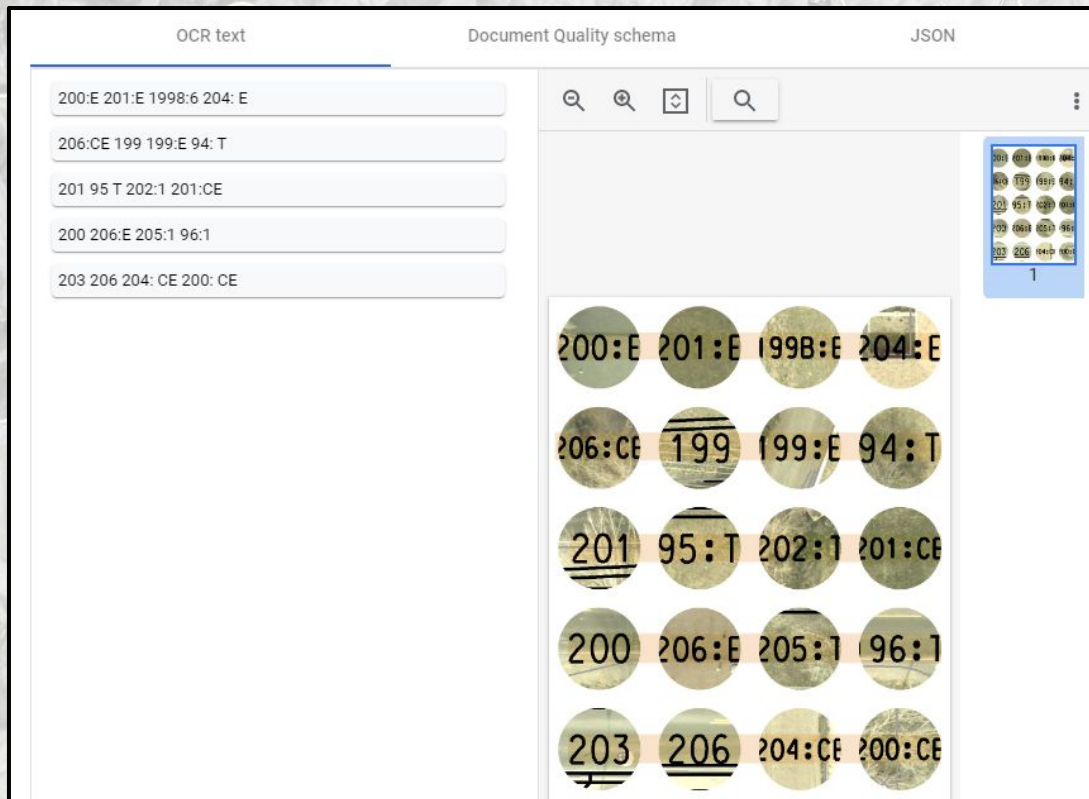


Crop & Mask
(original image)



Step 3 - Detect Text

- [Google Document AI](#) tool used to perform Optical Character Recognition (OCR) on each mosaic file
 - API extracts text and other metrics, returns a JSON file with results
- Basic cleanup on JSON text string
 - Remove newlines, whitespace, empty text results
- Insert results into a dataframe, with filename
- Save dataframe as a CSV



The screenshot displays the Google Document AI OCR interface. On the left, under the 'OCR text' tab, a list of detected text strings is shown:

- 200:E 201:E 1998:6 204: E
- 206:CE 199 199:E 94: T
- 201 95 T 202:1 201:CE
- 200 206:E 205:1 96:1
- 203 206 204: CE 200: CE

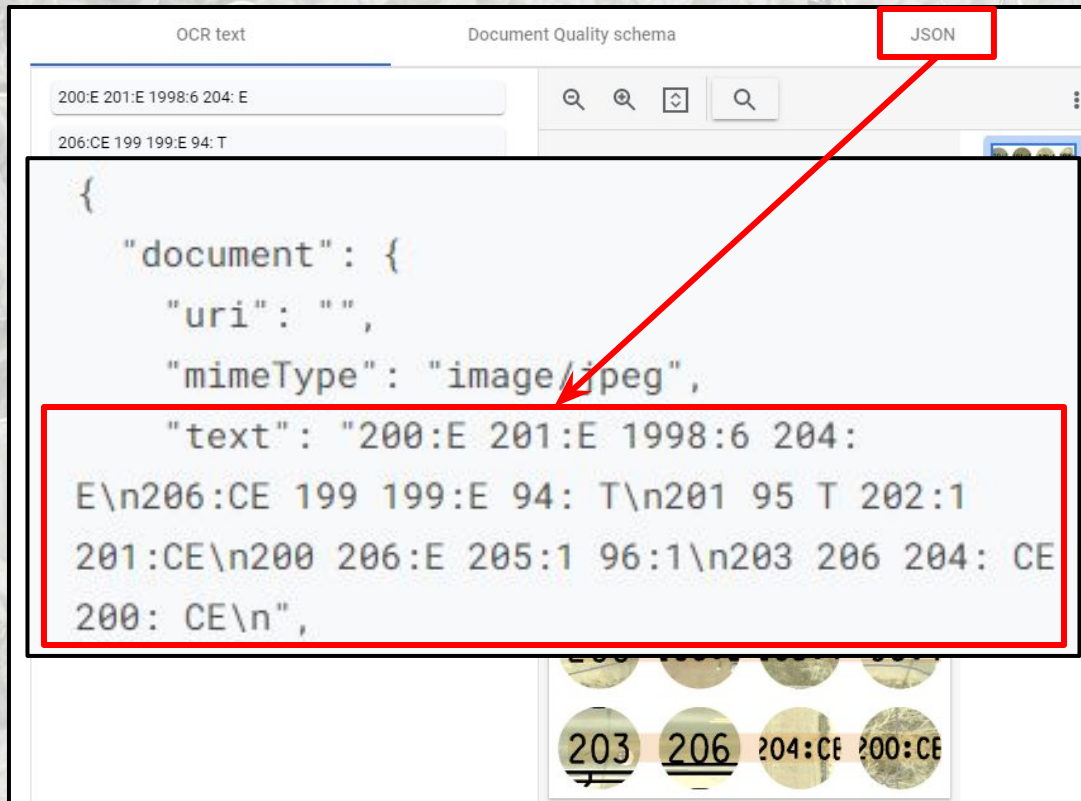
On the right, under the 'Document Quality schema' tab, a grid of circular image thumbnails is displayed. Each thumbnail contains a text overlay corresponding to the strings in the list. The text overlays are:

- 200:E 201:E 1998:E 204:E
- 206:CE 199 199:E 94: T
- 201 95: T 202:1 201:CE
- 200 206:E 205:1 96:1
- 203 206 204:CE 200:CE

 A small thumbnail grid on the far right shows a grid of 20 small circular images, with the first one highlighted and labeled '1'.

Step 3 - Detect Text

- [Google Document AI](#) tool used to perform Optical Character Recognition (OCR) on each mosaic file
 - API extracts text and other metrics, returns a JSON file with results
- Basic cleanup on JSON text string
 - Remove newlines, whitespace, empty text results
- Insert results into a dataframe, with filename
- Save dataframe as a CSV



The screenshot shows the Google Document AI OCR results interface. At the top, there are tabs for 'OCR text' and 'Document Quality schema', with a 'JSON' tab highlighted by a red box. Below the tabs, there are search and zoom controls. The main content area displays a JSON object representing the OCR results. A red arrow points from the 'JSON' tab to the JSON object. The JSON object is as follows:

```
{
  "document": {
    "uri": "",
    "mimeType": "image/jpeg",
    "text": "200:E 201:E 1998:6 204:
E\n206:CE 199 199:E 94: T\n201 95 T 202:1
201:CE\n200 206:E 205:1 96:1\n203 206 204: CE
200: CE\n",
```

Below the JSON object, there is a preview of the document image with text boxes overlaid on it. The text boxes contain the following text: 203, 206, 204:CE, 200:CE.

Step 4 - Combine Results

- Combine all result dataframes into a single dataframe (concatenate)
- Join additional fields from UDOT spreadsheets on filename field
 - project_number
 - project_name
 - guid from project management system (ProjectWise)
- Build URLs to ProjectWise, Cloud Storage files
- Explode text into multiple rows

	F	G	H	I	J	K	L	M
1	udot_file_name	project_number	project_name	guid	projectwise_url	udot_url	mosaic_url	text
2	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-proj	https://storage.	https://storag	21:4
3	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	24:E
4	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	25
5	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	24
6	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	25:3E
7	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	25:E
8	archive1/SP-0089(88)313_RW-15_Plan_Rev-1.TIF	SP-0089(88)313	SR-89, State Street,	{c88d59d5-f1a6-442e-a6b3-0bdebbdf3c88}	https://connect-pro	https://storage.	https://storag	24:3E
9	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	25:E
10	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	24:3E
11	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	24:E
12	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	25
13	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	25:3E
14	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	24
15	archive1/SP-0089(88)313_RW-15_Plan_Rev-2.TIF	SP-0089(88)313	SR-89, State Street,	{7f8fab8f-947e-4d06-aa4a-71f1505c43dd}	https://connect-proj	https://storage.	https://storag	21:4
16	archive1/SP-0089(88)313_RW-15_Plan_Rev-3.TIF	SP-0089(88)313	SR-89, State Street,	{eb9fc3b1-1d03-432f-800c-a39e0411b310}	https://connect-pro	https://storage.	https://storag	25

Step 5 - Post-process Results

- **Perform additional data cleanup**
 - Upper-case all letters
 - Remove punctuation (except for colons)
 - Remove character accents
- **Apply filtering rules to remove "invalid parcels", flag results:**
 - with special characters
 - starting with a letter or non-digit
 - with ':P' pattern
 - with a colon, if a number is not present before the colon
 - longer than 13 characters
 - with 4 or more letters, if no colon is present
 - with 5 or more numbers in a row
- **Remove duplicate rows**
- **Export final results into "good", "bad", and "all" spreadsheets**



UDOT Project Tools

- **OpenCV**

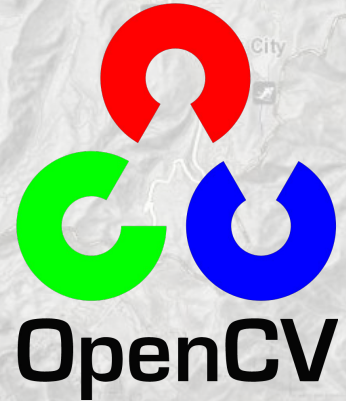
- Image manipulation
- Circle detection
- Cropping
- Mosaicking

- **Optical Character Recognition (OCR)**

- PyTesseract (original choice)
- Google DocumentAI (better results)

- **Pandas**

- Tabular data manipulation
- Joins
- String cleanup
- Filtering rules



github.com/agrc/udot-parcel-ml/

gis.utah.gov

UDOT Results

- Over **240,000 "good" parcels** were extracted from the documents

- **Data Accuracy**

- Reviewed 50 documents to compare CV results to human results
- Two "outlier" PDFs with 44 and 48 pages - might skew results

- **Circle detection**

- **Including outliers:** CV detected 1039/1244 circles (**83.52%**)
- **Excluding outliers:** CV detected 662/698 circles (**94.84%**)

- **Text comparison on 897 valid parcels**

- Average edit distance: 0.229
- Average **correct letter percentage** $(\text{truth_len} - \text{edit_dist}) / \text{truth_len}$: **95.15%**
- Number of results that were **perfect**: 878 (**88.96%**)
- Number of results with **edit distance ≤ 1** : 921 (**93.31%**)



Motivation - DHHS Cooling Towers

- Legionella bacteria can cause a serious type of pneumonia called Legionnaires' disease
- Legionella can grow/spread in large building water systems
 - Water tanks, HVAC, large/complex plumbing systems, cooling towers
- Cooling towers are concerning because they can release aerosolized water into the atmosphere
 - If Legionella is present, the aerosolized water can spread the bacteria over miles*

*[CDC - Controlling Legionella in Cooling Towers](#)

Commons Sources of Infection

Outbreaks of Legionnaires' disease are often associated with large or complex water systems, like those found in hospitals, hotels, and cruise ships.

The most likely sources of infection include:



Water used for showering (potable water)



Cooling towers (parts of large air conditioning systems)



Decorative fountains



Hot tubs

[CDC - Legionnaires' Fact Sheet](#)

Motivation - DHHS Cooling Towers

- Cooling towers can cause outbreaks of Legionnaires' disease when they are not adequately maintained
- They are often investigated & located using aerial imagery during Legionnaires' outbreaks
- Cooling towers have distinctive features that make them identifiable
- Researchers and the CDC have used object-detection models to identify potential cooling towers in aerial imagery ([TowerScout](#))



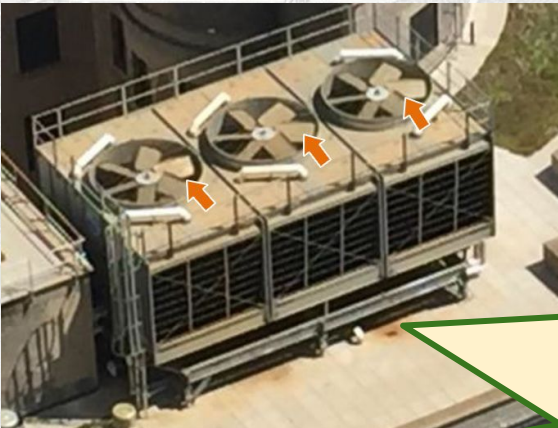
[CDC - Photos of Cooling Towers](#)

Motivation - DHHS Cooling Towers

- Cooling towers can cause outbreaks of Legionnaires' disease when they are not adequately maintained

This problem is solvable with computer vision!

and the CDC have used object-detection models to identify potential cooling towers in aerial imagery



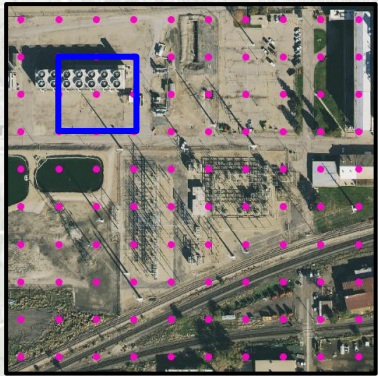
[CDC - Photos of Cooling Towers](#)

DHHS Cooling Towers - Process Overview

- 1) Build Index & Footprint
- 2) Download Images
- 3) Detect & Locate Towers
- 4) Post-process results
- 5) Build Web Map



1



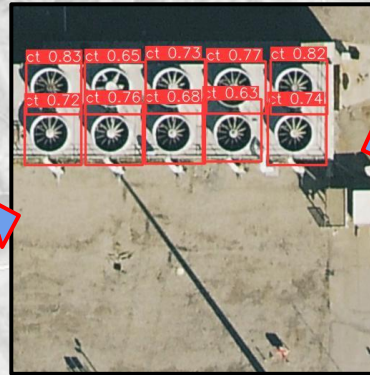
Build Index

2



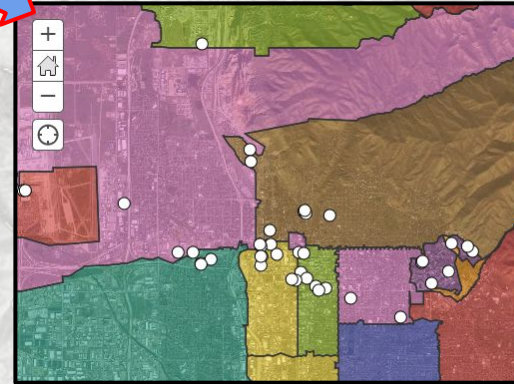
Download
Images

3



Detect &
Locate
Towers

4/5



Build Web
Map

Step 1 - Build Imagery Index & Footprint

- **Discover** Web Map Tile Services (WMTS) used for statewide imagery
- **Must build imagery index at highest zoom level (20)**
- **Create point for upper left corner of each WMTS tile (row, col)**
 - Calculate latitude/longitude based on row/col values
 - 275,000,000 tiles cover the state of Utah!
- **Create processing footprint to select a subset of all image tiles**
 - Buffer census places by 800m
 - Buffer large buildings (>5k sq ft) by 800m
- **Select WMTS indices of tiles in processing footprint**
 - Load tile index and footprint into Google BigQuery
 - Run Spatial SQL query to select tiles to process (~6%)
 - (data is waaaayyy too big for desktop GIS)



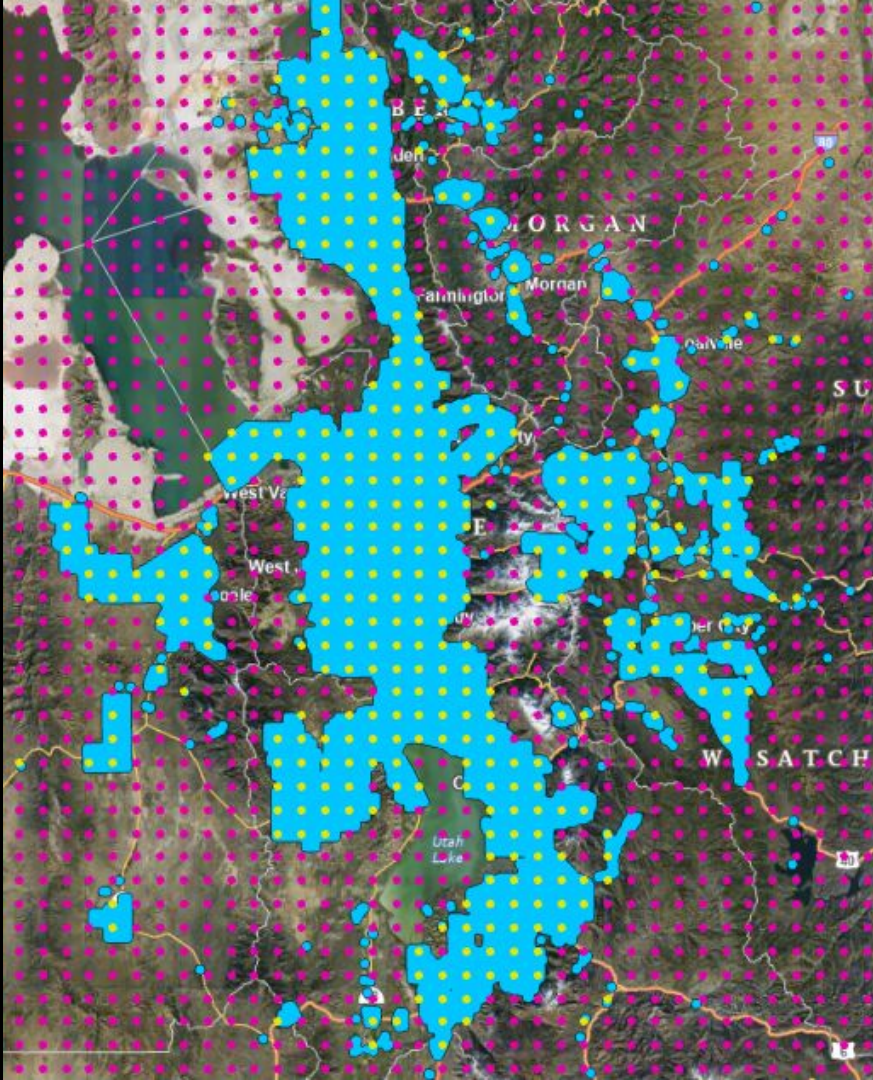
Build Index



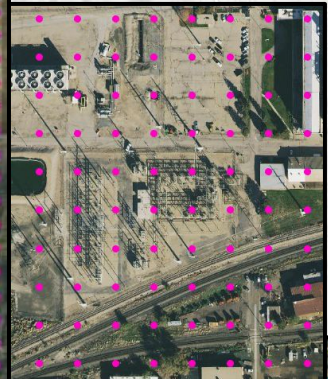
Create Footprint

Step 1 - Bu

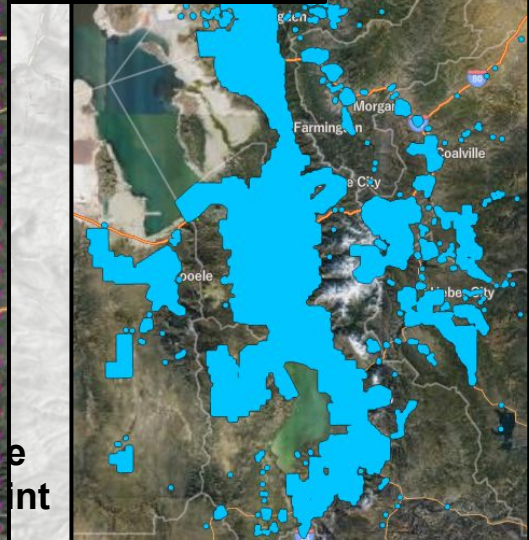
- [Discover](#) Web Ma
statewide imagery
- Must build image
- Create point for u
(row, col)
 - Calculate latitude
 - 275,000,000 tiles
- Create processing
image tiles
 - Buffer census pla
 - Buffer large build
- Select WMTS indi
 - Load tile index an
 - Run Spatial SQL
 - (data is waaaayy)



Footprint



Build
Index



e
int

Step 2 - Download Images



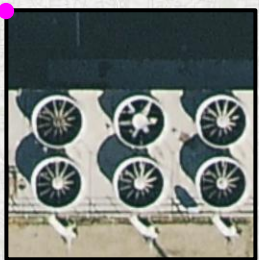
- Iterate through tile indices within footprint
 - Download primary tile and 3 neighboring tiles with HTTPS GET requests:
 - https://discover.agrc.utah.gov/login/path/{quad-word}/tiles/15cm_hexagon_utah/20/{col}/{row}

- Build mosaic image

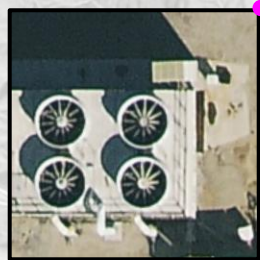
Each WMTS tile is 256x256 pixels

Model input performs best on 512x512 images

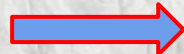
{col}, {row}



198263, 394029



198264, 394029



198263, 394030

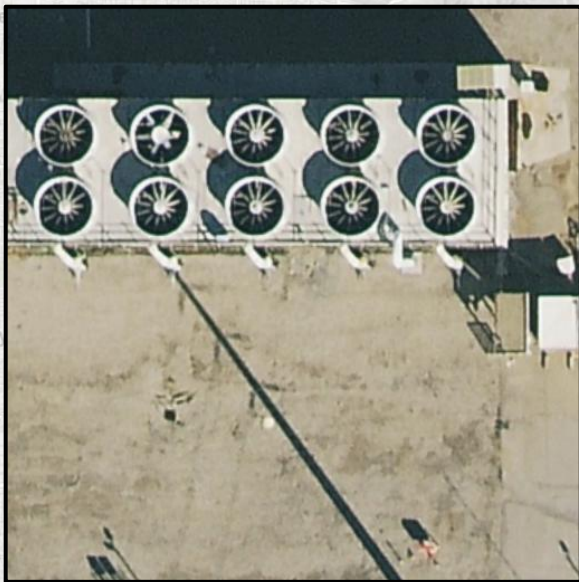


198264, 394030

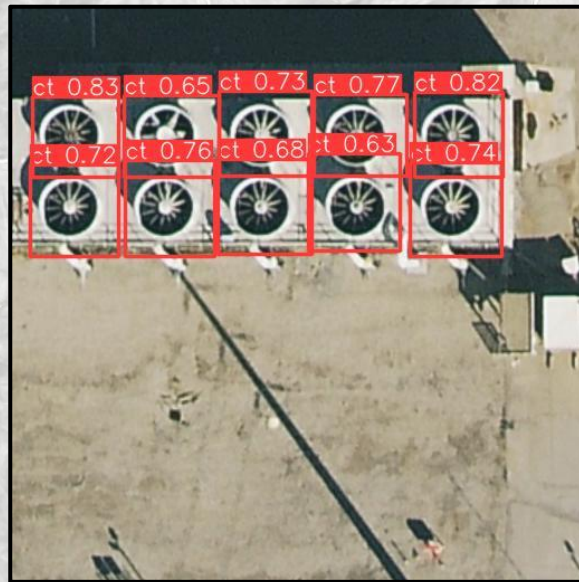
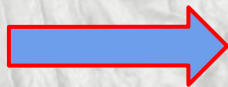


Step 3 - Detect & Locate Towers

- Run PyTorch model on each mosaic
 - Pre-trained "[TowerScout](#)" model, provided by CDC
 - YOLOv5 backbone
- Get results as a dataframe



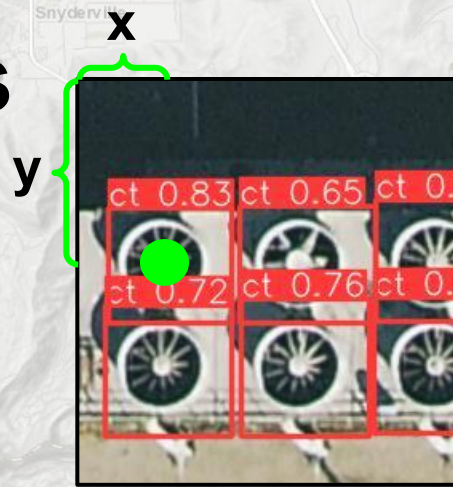
Detect
Towers



*TowerScout created by Karen Wong, Jia Lu, Gunnar Mein, and Thaddeus Segura, licensed under [CC-BY-NC-SA-4.0](#)

Step 3 - Detect & Locate Towers

- Calculate X/Y coordinates of tower centroid
 - PyTorch model provides tower location bounding box in pixels from upper-left corner
 - Centroid is calculated from xmin, xmax, ymin, ymax values
 - Able to convert pixels into geographic space because we know:
 - 1) coordinates of the tile's upper-left corner
 - 2) pixel size at WMTS zoom level 20 (0.1492910708688 meters)



	xmin	ymin	xmax	ymax	confidence	x_centroid_px	y_centroid_px	x_centroid_3857	y_centroid_3857
9	265.932	130.687	347.821	217.394	0.633	306.877	174.041	-12460145.259	4978279.451
8	100.593	78.307	185.93	151.344	0.654	143.262	114.826	-12460169.686	4978288.291
7	184.386	135.914	268.444	220.131	0.676	226.415	178.023	-12460157.272	4978278.857

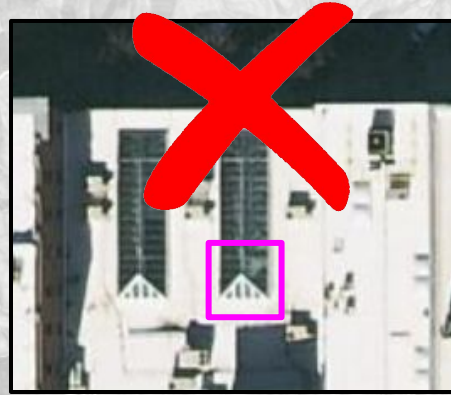
```
#: x, y = upper-left corner coordinates of tile in web mercator
#: calculate centroid x/y coords in web mercator
meters_per_pixel = 0.1492910708688
results_df["x_centroid_3857"] = x + results_df["x_centroid_px"] * meters_per_pixel
results_df["y_centroid_3857"] = y - results_df["y_centroid_px"] * meters_per_pixel
```

Calculate X/Y

- Upload results to BigQuery

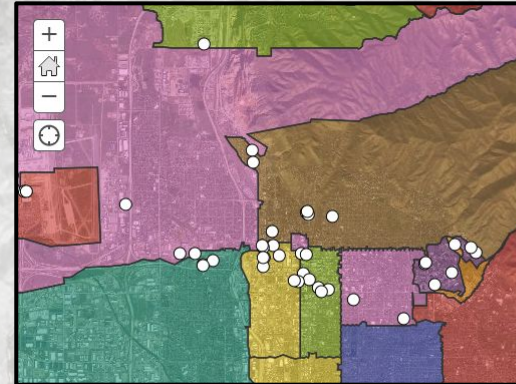
Step 4 - Post-process Results

- Manually validate detected cooling towers
- Enrich data with attributes from other statewide datasets
 - Nearby address
 - County
 - City
 - Zip Code
 - Small Health Statistical Area



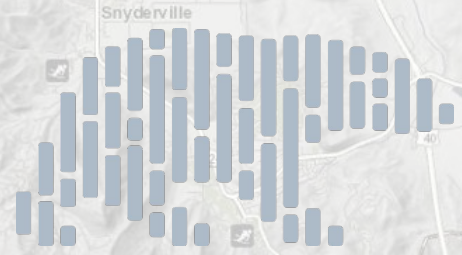
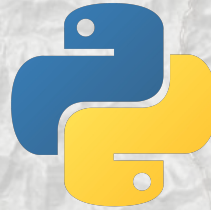
Step 5 - Build Web Map

- Create hosted feature layer of cooling tower locations
- Build a web map for DHHS users
- Add other relevant, health-related layers and tools



DHHS Tools

- **Mercantile** - WMTS tile index to lat/lon
- **Polars** - large dataframe creation
- **Requests** - http requests and downloads
- **Pyproj** - coordinate conversion
- **PyTorch** - object detection model
- **Google Cloud Platform**
 - **BigQuery** - massive tabular data, querying, spatial SQL
 - **Cloud Run** - cloud computing



Requests
http for humans



github.com/agrc/dhhs-cooling-towers

gis.utah.gov

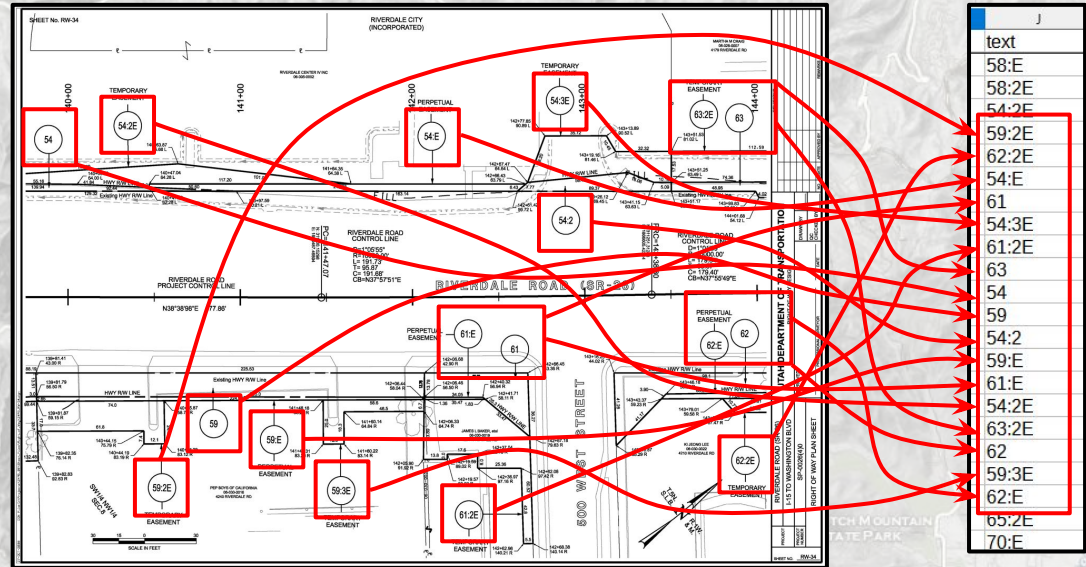
DHHS Project Results

- To be determined...
- Most of the code is written, but the processing is in-work
- Check back later for an update! (UGRC blog post)



Final words

- A lot of useful data and information can be locked away in documents and imagery
- Computer Vision tools can help unlock that data
- Cloud computing can reduce processing time by orders of magnitude
- UDOT parcel detection and DHHS cooling towers projects highlight these possibilities



Questions?



Location matters

Erik Neemann



email: eneemann@utah.gov



twitter: [@Erik_UGRC](https://twitter.com/@Erik_UGRC)