

# Character Recognition in < 100 Py Statements

**Ash Richardson**  
**Senior Data Scientist**  
**Predictive Services Unit**

Self-contained implementation "no" dependencies



BC Wildfire  
Service

“The syntactical nature of reality, the real secret of magic, is that the world is made of words. And if you know the words that the world is made of, you can make of it whatever you wish.”

Terrence M (1946 – 2000) Writer, Ethnobotanist, Mystic, Philosopher (USA)

# Welcome to Python online and DS CoP!



- Everyone welcome
- Thank you for coming!
- Offer to do a talk or demo Please?
- Door is always open.. Feedback / input anytime pls.!
- Check out the code together session! / Experts on tap
- Please join channels and participate in discussions + events?

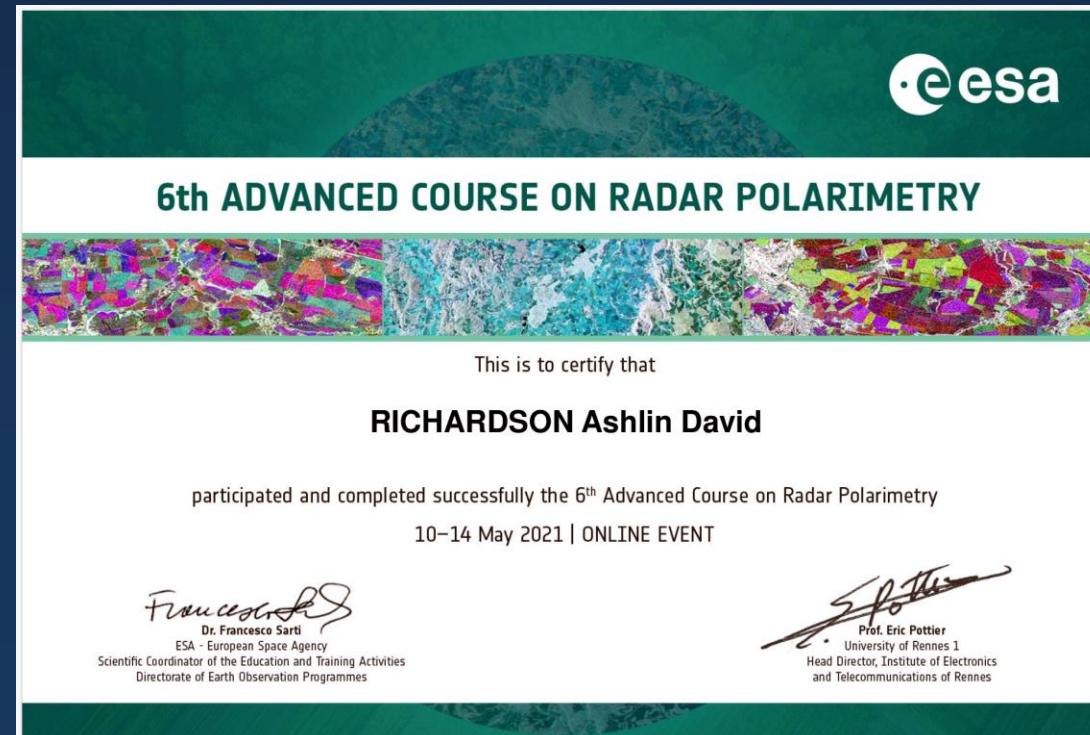
# Overview



- Explore Character Recognition -- a "core" Machine Learning task
  - Construct accessible implementation
  - Offer some baseline of ML algorithms complexity
    - "Keep it simple"
    - Flexible and scalable (to larger data)
    - Incremental springboard to support communicating other results later..
    - “base” Python with “no” dependencies
    - < 100 lines (ish)..
  - Stay as non-technical as possible
  - Keep Emphasis on discussion & relationships & connections..

# Plan

- Some context
  - Why explicit, self-contained?
    - Control
    - Predictable, reliable
    - Want progs to run for 10, 20, 30+ years?
- ML Approach     Dive in...
  - 1) Segment
    - "Flood-fill" segmentation
  - 2) Classify
    - Earth Mover Distance (EMD) (ish)..
- Visualize 1) and 2)
  - Iterate on some examples

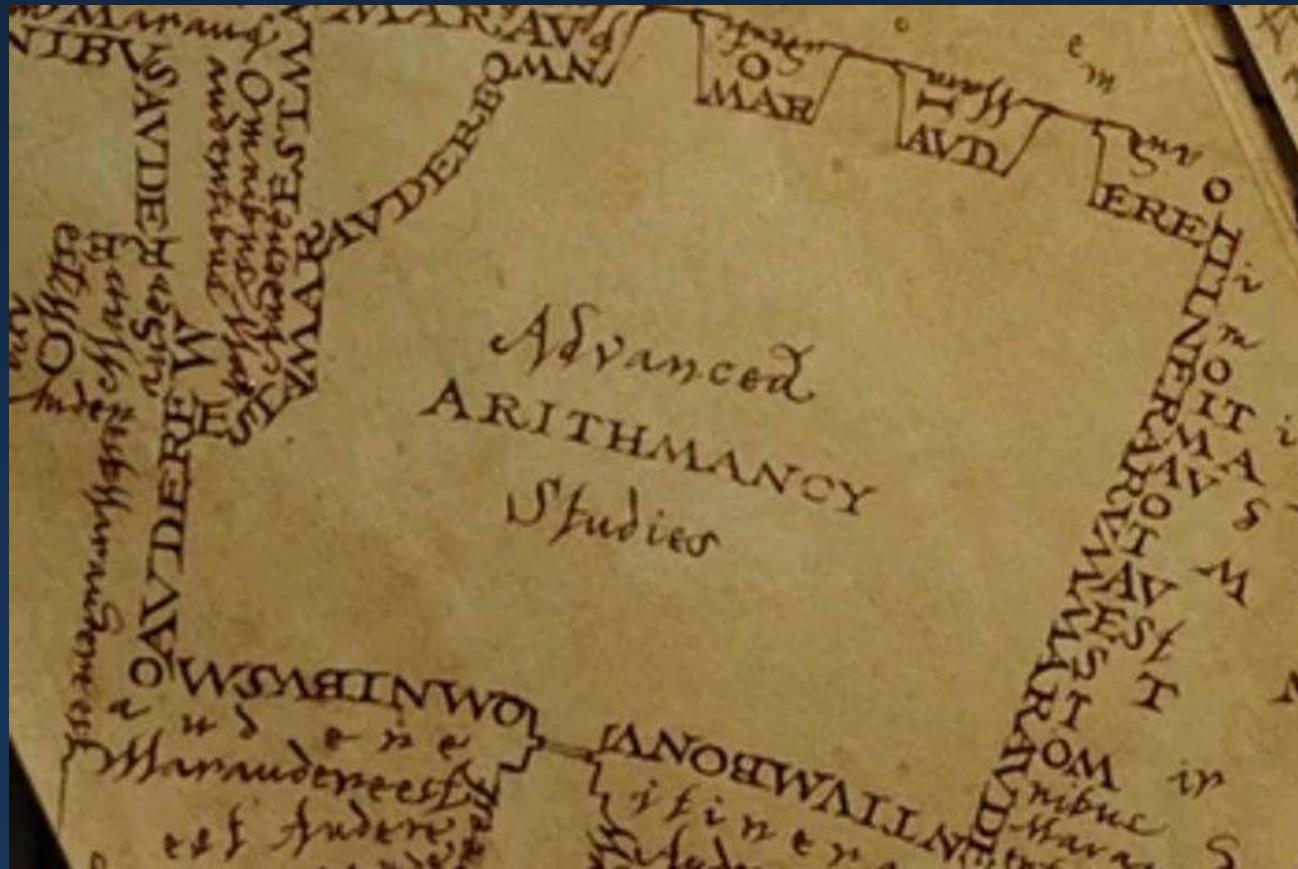


# What kind of magic?



BC Wildfire  
Service

Building maps with symbols..



(English... transport planning  
for  
statistical  
predictive analytics..)

Transfiguration  
for  
Arithmantic  
divination?

## CURRICULUM

### MANDATORY CLASSES

- TRANSFIGURATION
- DEFENCE AGAINST THE DARK ARTS
- CHARMS
- POTIONS
- ASTRONOMY
- HISTORY OF MAGIC
- HERBOLOGY
- FLYING (FIRST YEARS ONLY)

### ELECTIVE CLASSES [THIRD YEARS & UP CHOOSE TWO]

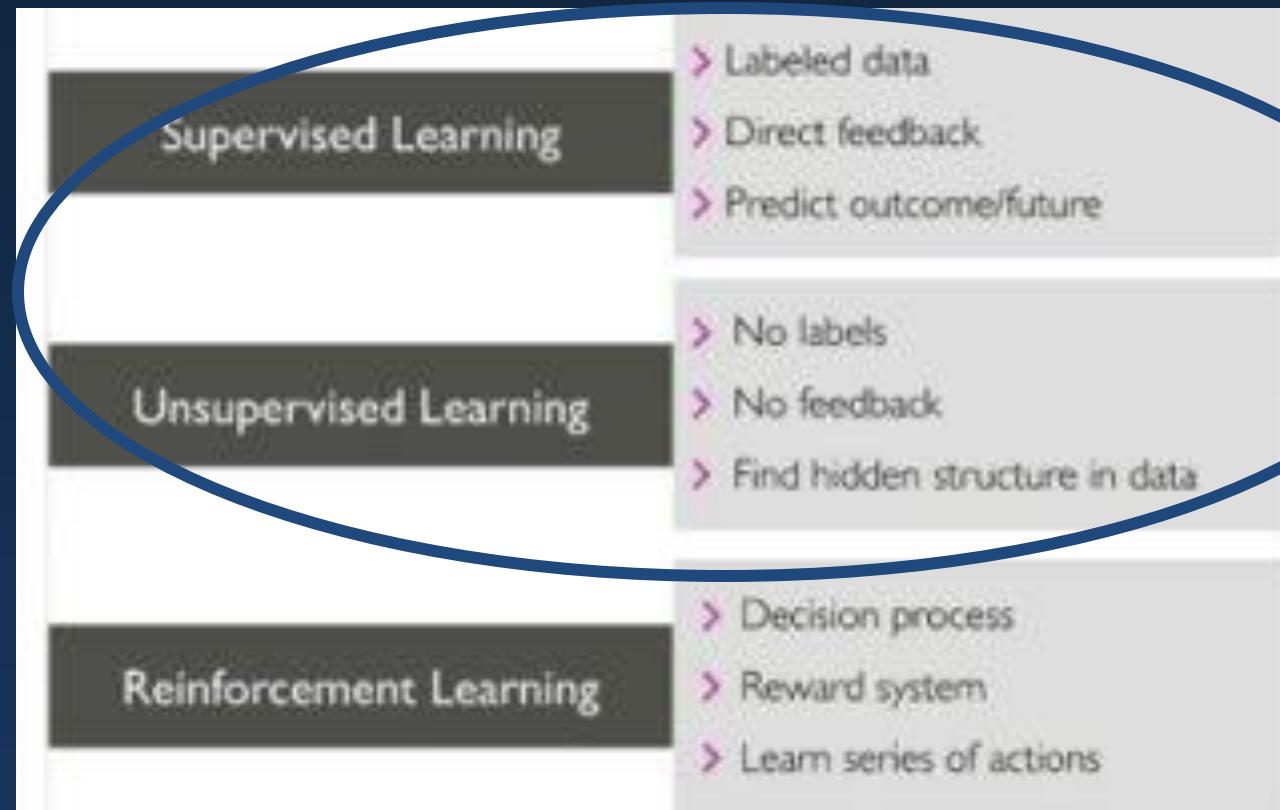
- ALCHEMY
- APPARITION
- ARITHMANTIC
- ANCIENT RUNES
- CARE OF MAGICAL CREATURES
- DIVINATION
- MUGGLE STUDIES
- MUSIC



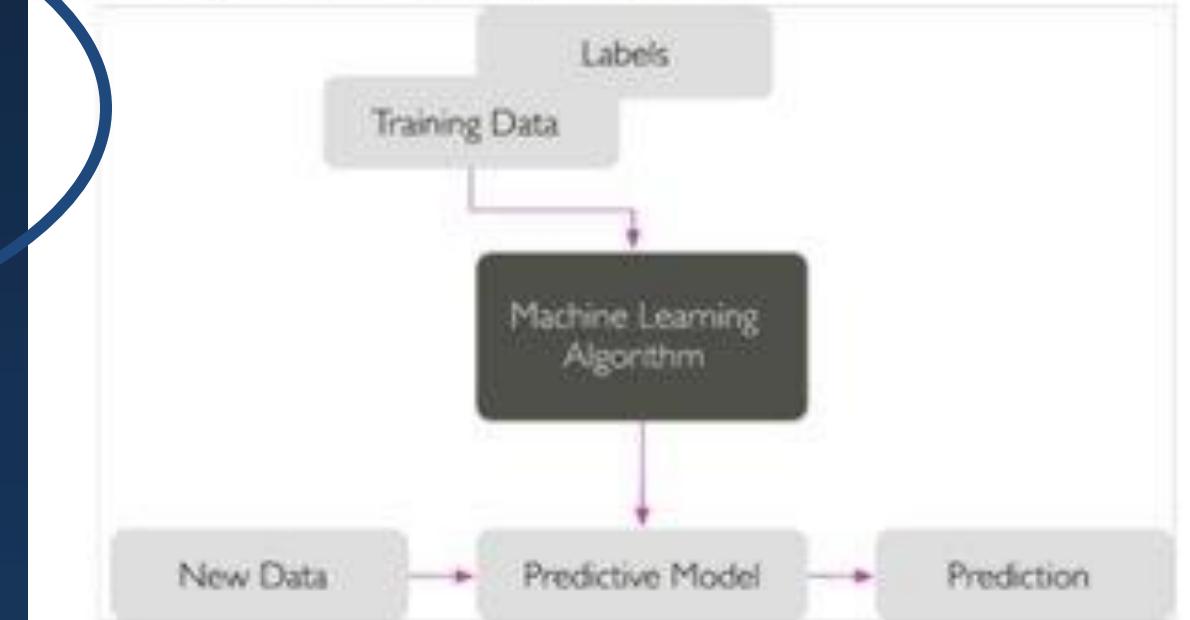
# Character recognition, what kind of ML?



BC Wildfire  
Service



Making predictions about the future with supervised learning  
The main goal in supervised learning is to learn a model from labeled training data that allows us to make predictions about unseen or future data. Here, the term **supervised** refers to a set of samples where the desired output signals (labels) are already known.



Supervised:

Get samples of labelled data,  
try to label new data and get the right answer..

From "KD Nuggets" blog

# Flow chart

- 0) Rendering
  - Generate "raw binary" truth & test data.. LaTeX disable ligatures
- 1) Segmentation (flood fill)
  - Cut up truth & test images (read them with no driver)
- 2) Predict ("Earth mover" style distance)
  - Compare test data to training data
    - Find "nearest" truth data

## render.py

```
1 ''' render some text in latex, converting the typeset result to "raw binary"
2 ..which can be read without a driver'''
3 import os
4
5 def run(c): # run something at terminal and wait to finish
6     print(c)
7     a = os.system(c)
8
9 def render(my_text, name):
10    run('mkdir -p ' + name)
11
12    # use "Computer Modern" font by Donald Knuth # insert Knuth quotes..
13    open(name + '.tex', 'wb').write('\n'.join(['\\documentclass{letter}', 
14        '\\usepackage{lmodern}', 
15        '\\usepackage[T1]{fontenc}', 
16        '\\usepackage{xcolor}', 
17        '\\usepackage[tracking=true]{microtype}', 
18        '\\DisableLigatures{encoding=any}', 
19        '\\pagenumbering{gobble}', 
20        '\\begin{document}', 
21        '\\color{blue}'] + # blue
22        [ '\\textsf{' + my_text +
23
24    if not os.path.exists(name + '.bin'): # delete train.bin to start from scratch
25        run('pdflatex ' + name + '.tex') # render with LaTeX
26        run('convert -background white -density 333 ' + name + '.pdf' + name + '.bin')
27        run('gdal_translate -of ENVI -ot Float32 ' + name + '.bmp' + name + '.bin')
28
29    if os.path.exists(name + '.hdr'):
30        d = open(name + '.hdr').read() + 'band names = {red,\ngreen,\nblue}'
31        open(name + '.hdr', 'wb').write(d.encode())
32
33 truth = [] # truth characters
34
35 if __name__ == "__main__":
36
37    def chars(i, j): # add chars between ascii codes i, j to truth data
38        global truth
39        my_chars = [chr(x) for x in range(i, j)]
40        truth += my_chars
41        print("my_chars", my_chars)
42        return ''.join(my_chars)
43
44    if not os.path.exists('truth.bin'):
45        render([chars(48, 58) + '\n', chars(65, 91) + '\n', chars(97, 123) + '\n'])
46        'truth' # designate as truth
47
48        print('+w truth_chars.txt')
49        open('truth_chars.txt', 'wb').write(''.join(truth).encode()) #
50
51    if not os.path.exists('test.bin'):
52        print("render test data..")
53        # render(["hello world"], 'test')
54
55    render(["Through three cheese trees\\\\ \\\\", 
56            "three free fleas flew\\\\ \\\\", 
57            "While these fleas flew\\\\ \\\\", 
58            "freezy breeze blew\\\\ \\\\", 
59            "Freezy breeze made\\\\ \\\\", 
60            "these three trees freeze\\\\ \\\\", 
61            "Freezy trees made\\\\ \\\\", 
62            "these trees cheese freeze\\\\ \\\\", 
63            "Thats what made these\\\\ \\\\", 
64            "three free fleas sneeze\\\\ \\\\"], 
65            'test')
```

$AE \rightarrow \mathcal{A}\mathcal{E}$	$ij \rightarrow ij$
$ae \rightarrow \alpha$	$st \rightarrow st$
$OE \rightarrow \mathcal{O}\mathcal{E}$	$ft \rightarrow ft$
$oe \rightarrow \alpha$	$et \rightarrow \&$
$ff \rightarrow ff$	$fs \rightarrow \beta$
$fi \rightarrow fi$	$ffi \rightarrow ffi$

# 0) Rendering..

- 'pdflatex'
  - Ligatures off
  - Space the characters out slightly
- 'convert' (imagemagick)
  - Convert pdf to bmp
- 'gdal\_translate'
  - Convert bmp to "raw" binary

Result: beautiful images of text that  
Can be read without a driver..

# How to read image data without driver



```
def read_hdr(hdr): # read the image dimensions
    cols, rows, bands = 0, 0, 0
    for line in open(hdr).readlines():
        chunks = line.strip().split('=')
        try: # pull off two chunks delimited by '='
            f, g = [x.strip() for x in chunks[0:2]]
            if f == 'samples':
                cols = g
            if f == 'lines':
                rows = g
            if f == 'bands':
                bands = g
        except:
            pass
    return [int(x) for x in [cols, rows, bands]] # string to int

def read_float(fn): # read the raw binary file
    return np.fromfile(fn, dtype=np.float32) / 255. # put data in range [0, 1]

'''pixel @ (row, col) = (i, j):
npx = nrow * ncol # number of pixels in image
red value: dat[i * ncol + j]
grn value: dat[npx + i * ncol + j]
blu value: dat[2 * npx + i * ncol + j]'''
```

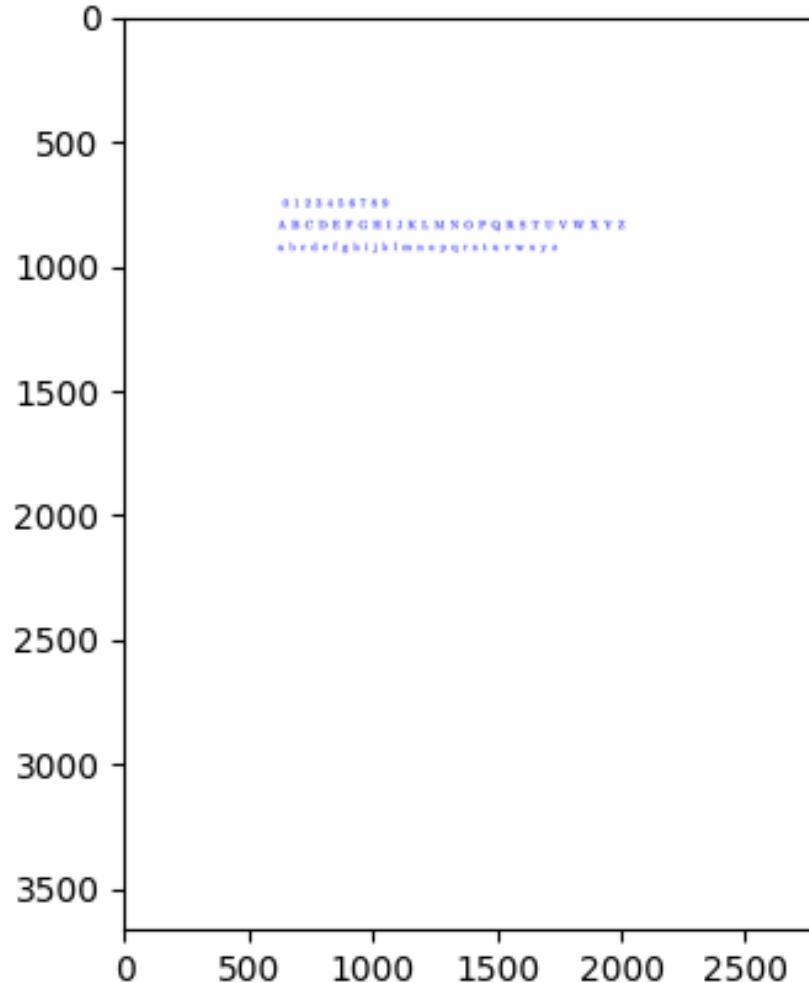
truth.hdr  
ENVI  
samples = 2831  
lines = 3663  
bands = 3  
header offset = 0  
file type = ENVI Standard  
data type = 4  
interleave = bsq  
byte order = 0  
band names = {red,  
green,  
blue}

"flat file" of image world? Lots  
of people prefer geotiffs..

# Truth



BC Wildfire  
Service

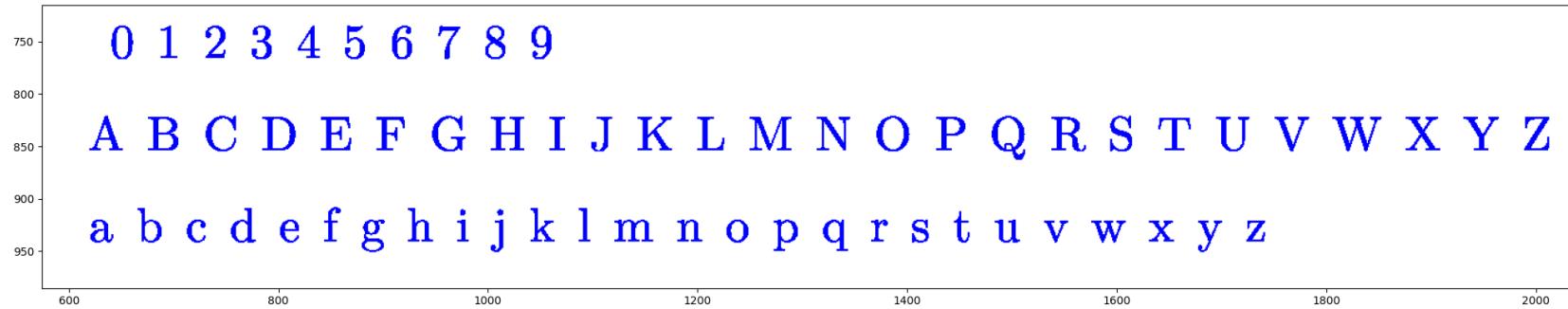


0123456789  
ABCDEFGHIJKLMNOPQRSTUVWXYZ  
abcdefghijklmnopqrstuvwxyz

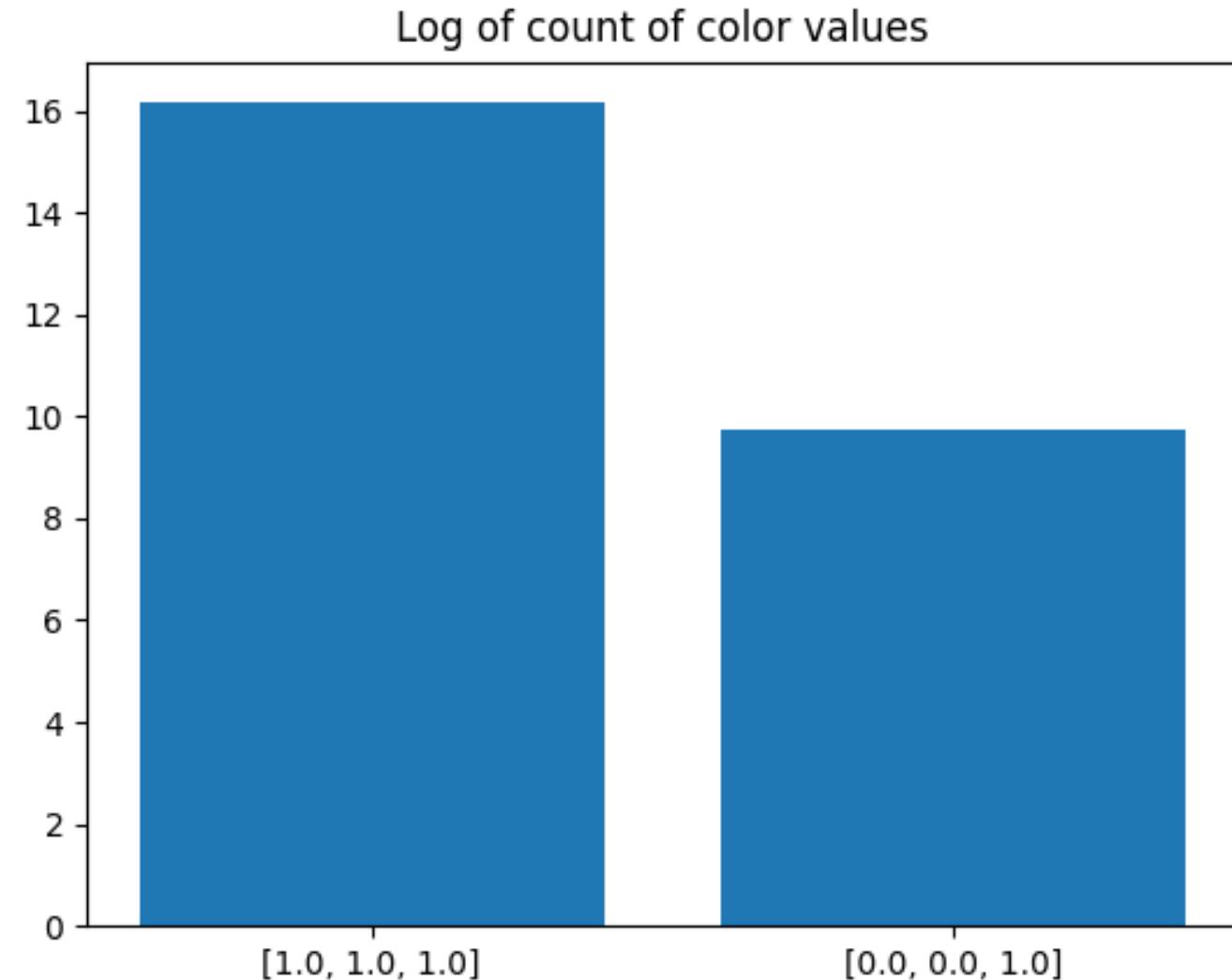
# Truth



BC Wildfire  
Service



# Truth



[1,1,1] is white  
[0,0,1] is blue

# Flood fill

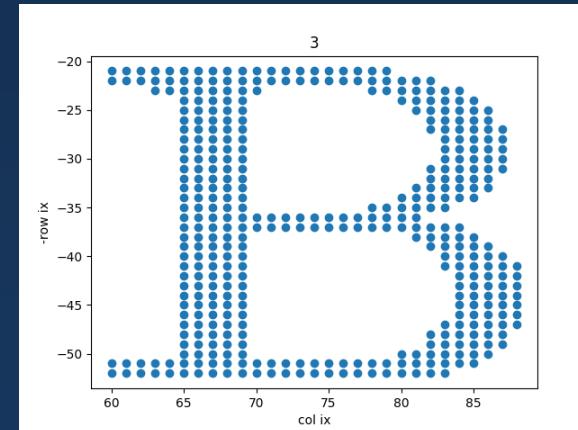
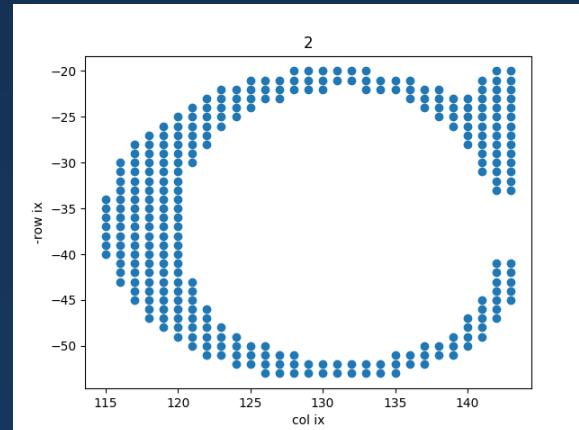
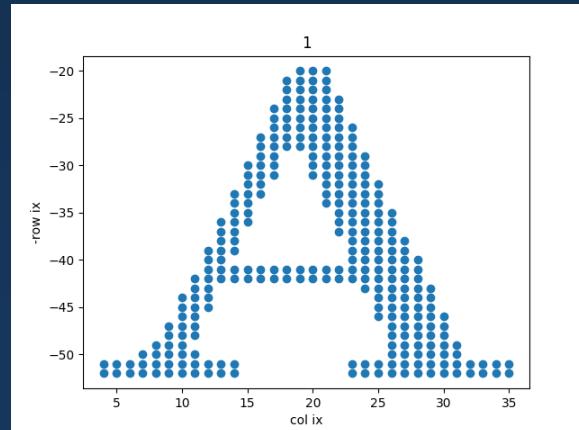


BC Wildfire  
Service

A      B      C

```
# Recap the last step (render.py does this for us)
gdal_translate -of ENVI -ot Float32 abc.bmp abc.bin
Input file size is 156, 75
0...10...20...30...40...50...60...70...80...90...100 - done.
```

```
$ python
>>> from image import image
>>> a = image("abc.bin")
>>> a.segment()
```



Point cloud  
representations

Why does B come last?

# Step 2. Flood fill...

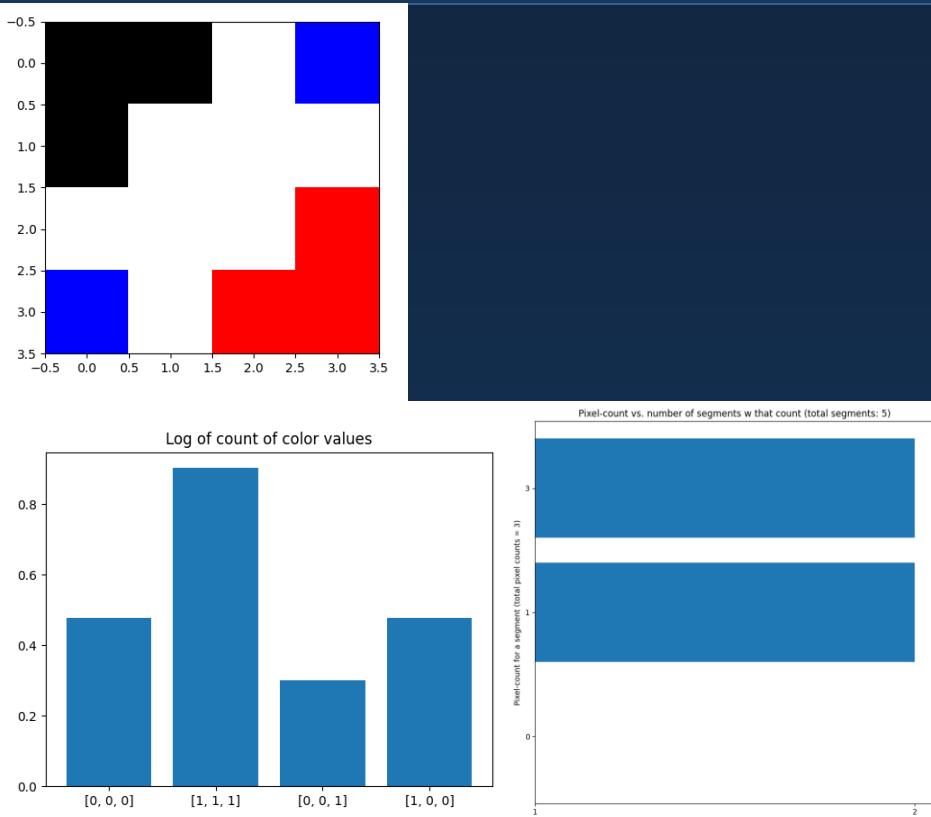


BC Wildfire  
Service

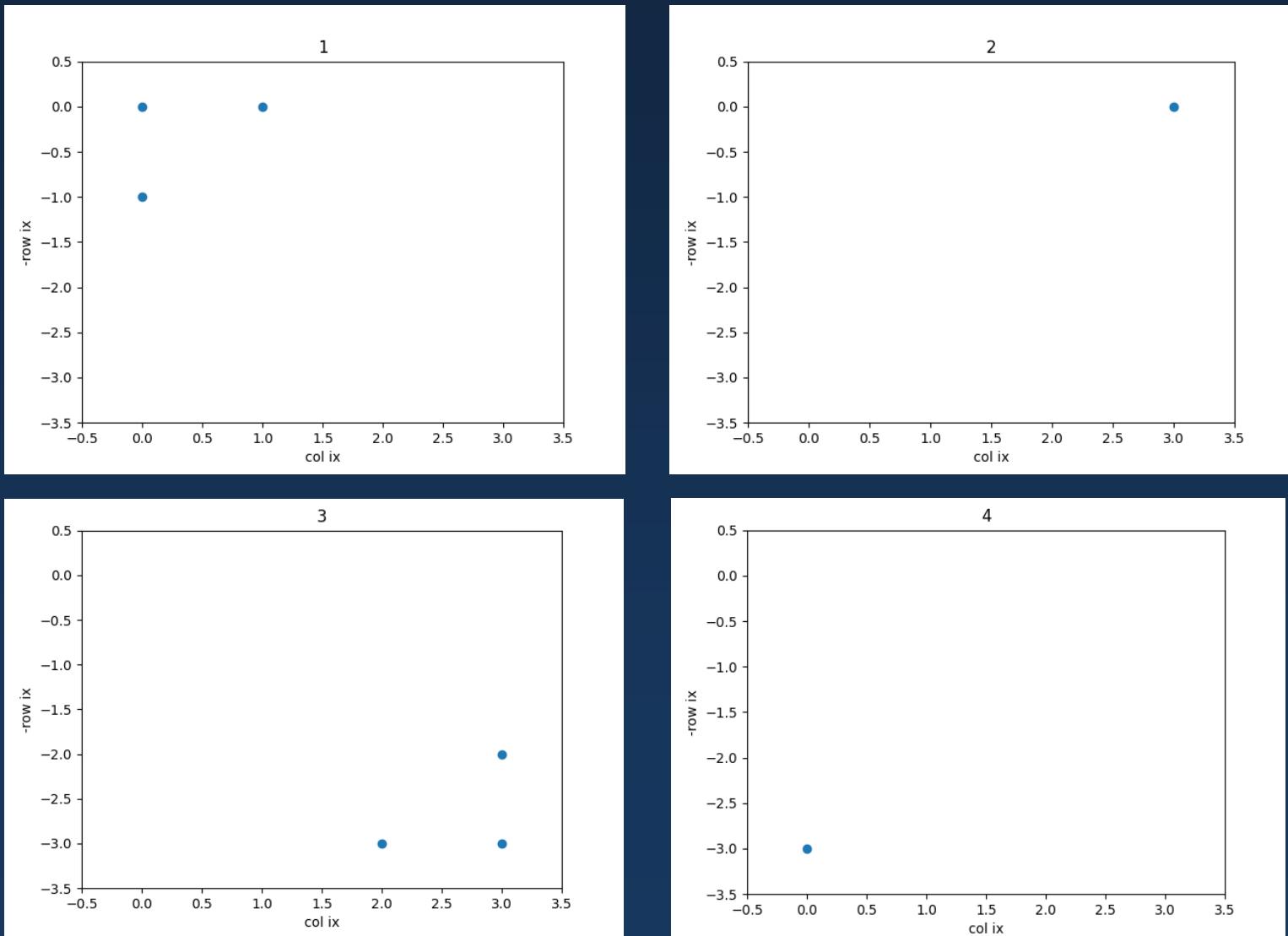
```
1 '''flood fill segmentation'''
2 def flood(img, i, j, my_label=None, my_color=None):
3     if i >= img.rows or j >= img.cols or i < 0 or j < 0:
4         return # out of bounds
5     ix = i * img.cols + j # linear index of (i, j)
6
7     if img.labels[ix] > 0:
8         return # stop: already labelled
9
10    c_s = str(img.rgb[ix])
11    if c_s == img.max_color:
12        return # stop: ignore background "background subtraction"
13    if my_color and my_color != c_s:
14        return # stop: different colour than at invocation chain start
15
16    if my_label:
17        img.labels[ix] = my_label
18    else:
19        img.labels[ix] = img.next_label
20        img.next_label += 1
21
22    for di in [-1, 0, 1]:
23        for dj in [-1, 0, 1]:
24            if not (di == 0 and dj == 0):
25                flood(img, i + di, j + dj, img.labels[ix], c_s)
```



# Step 2. Flood fill...



Why this order?



# Flood fill segmentation



BC Wildfire  
Service

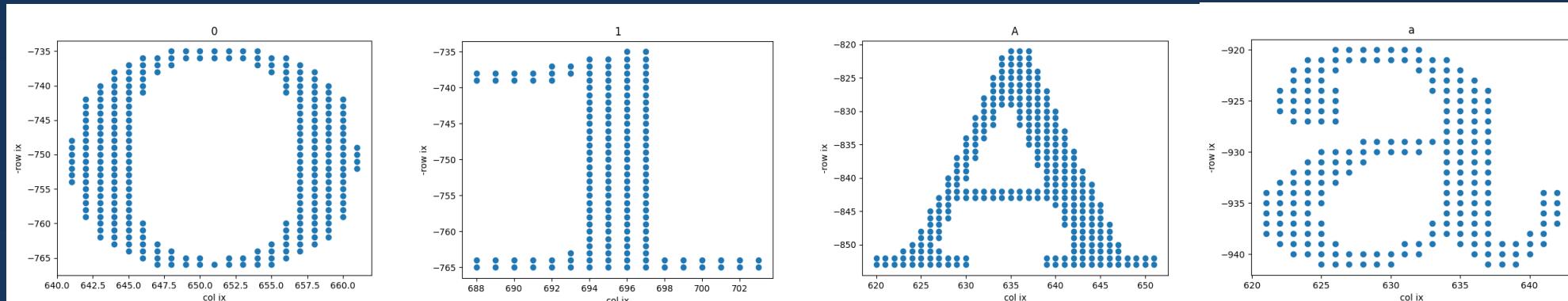
- "Smart scissors"
  - not "very" smart.. Auto scissors!
  - Grabs "Connected areas of uniform color"...
- Glyphs saved..
  - Truth/\*.p
  - Truth/\*.png (Check filenames match plots!)
  - Test/\*.p
  - Test/\*.png

segment.py

```
1 from image import image
2
3 # segment the truth data
4 truth = image('truth.bin')
5 truth.segment([745, 838, 932])
6
7 # segment the test data
8 test = image('test.bin')
9 test.segment()
```



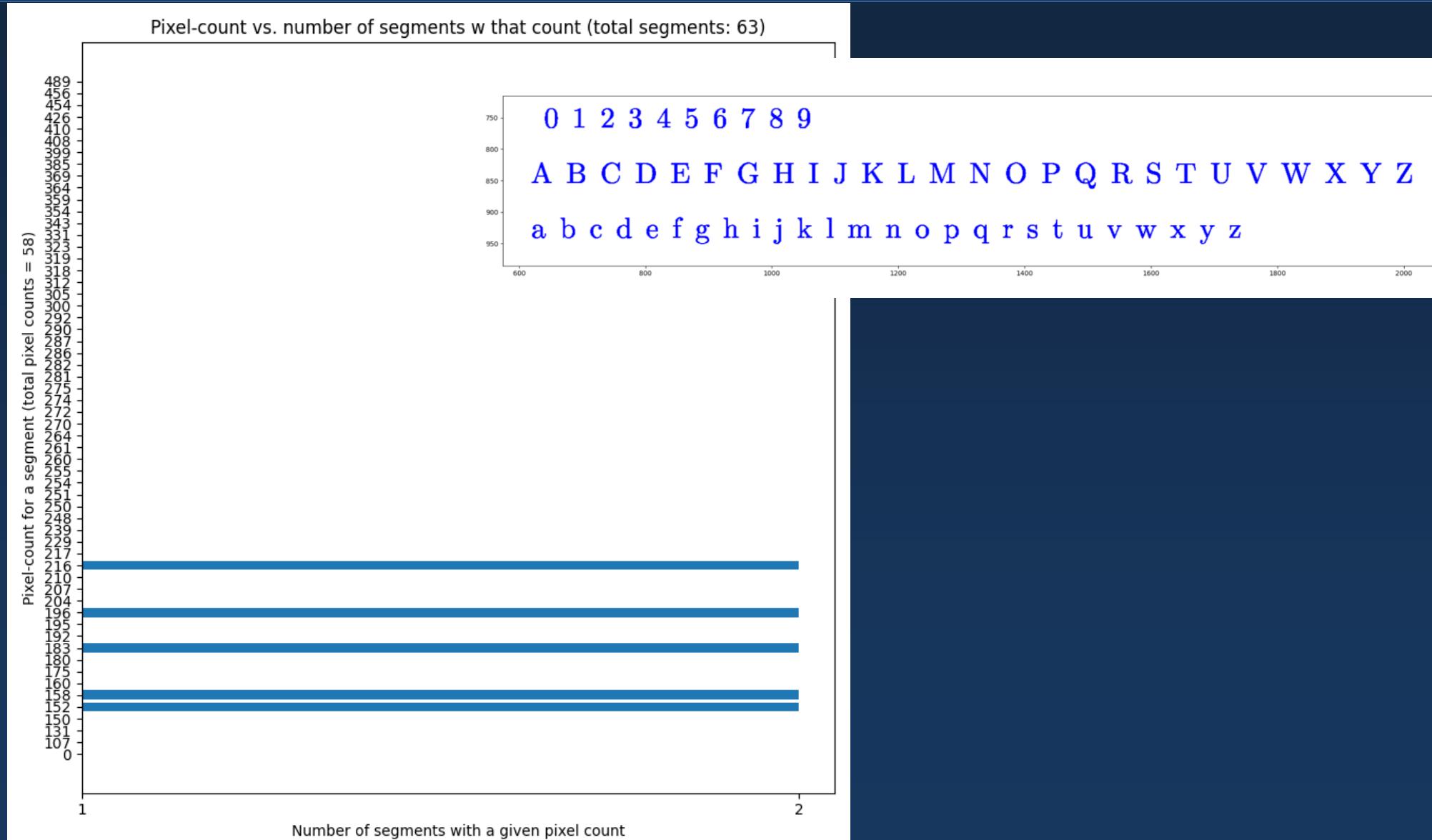
why?



# Truth

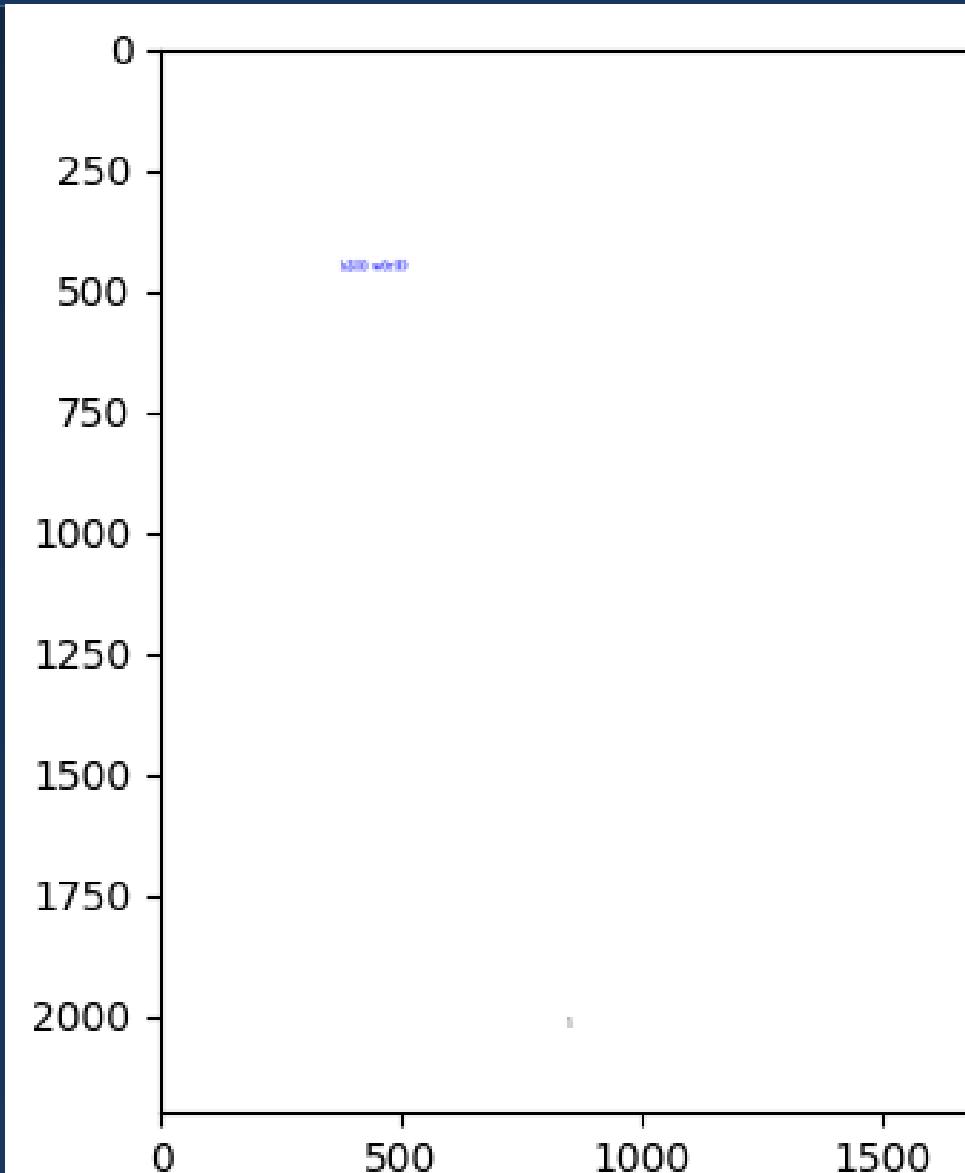


BC Wildfire  
Service

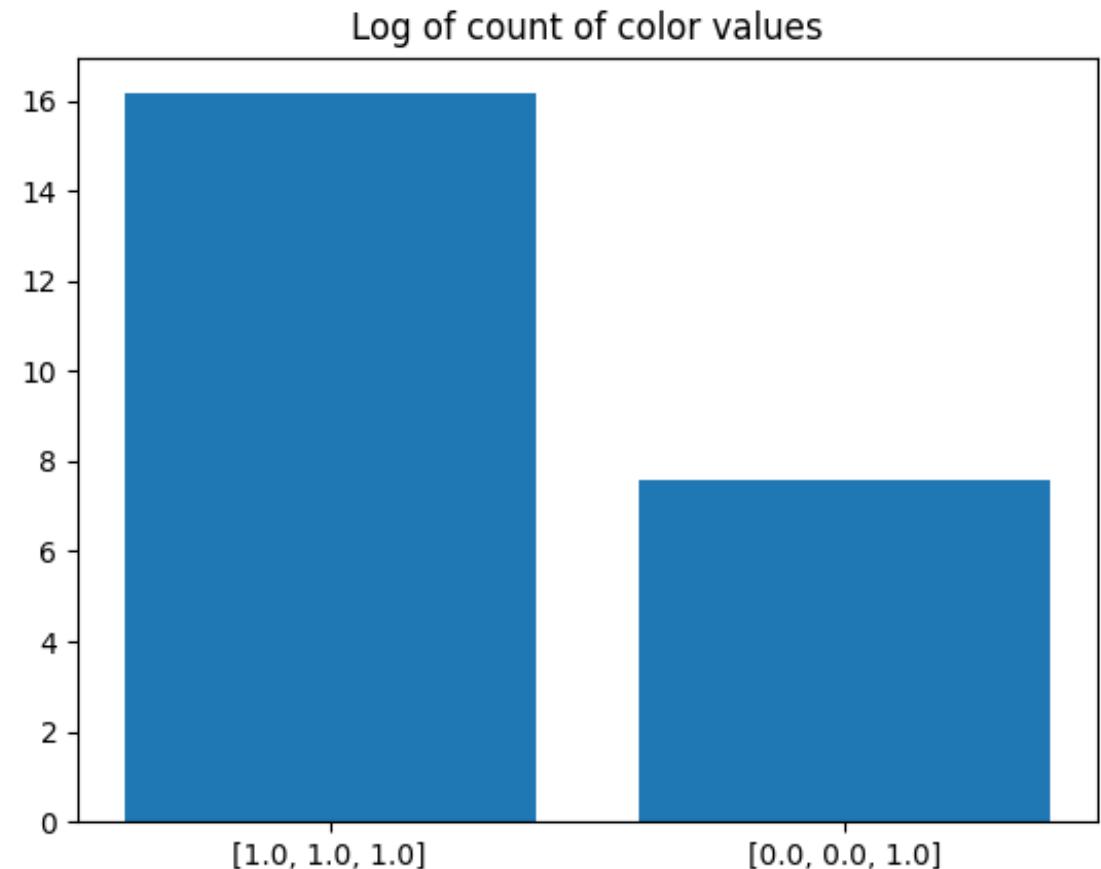
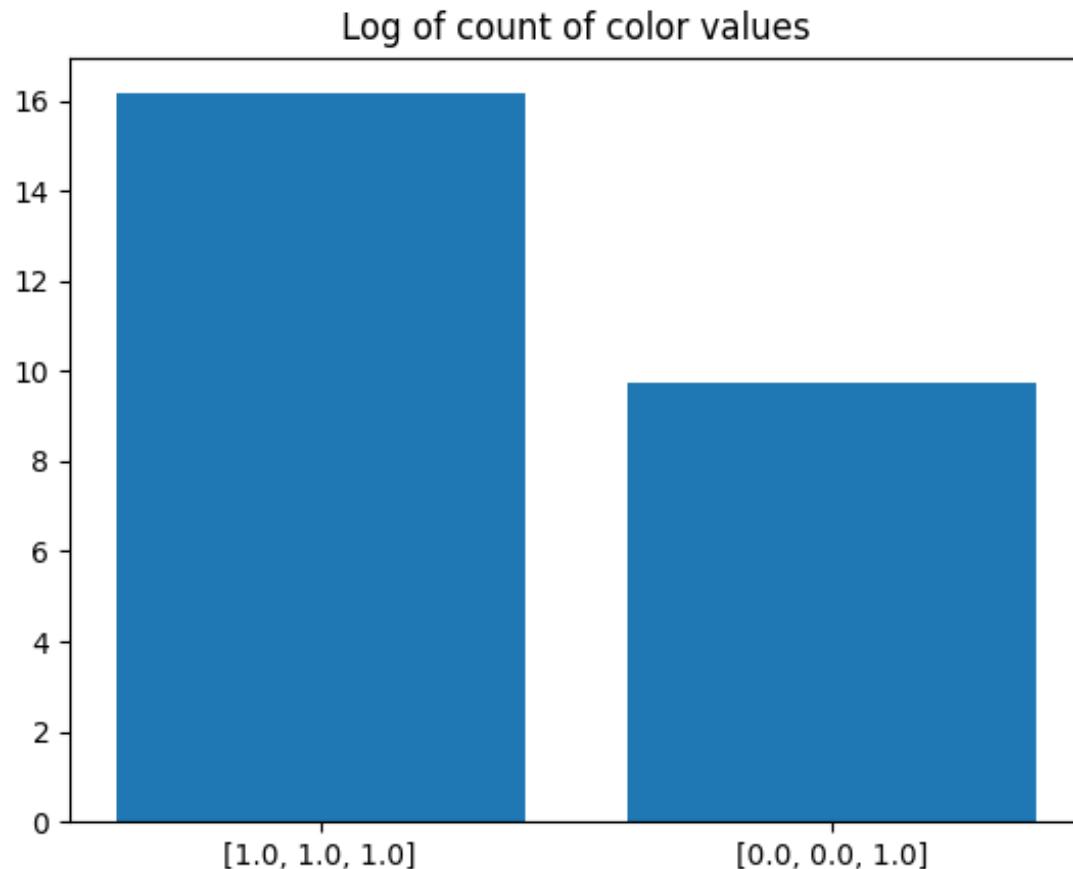


Most "segments" have  
Unique pixel-counts..  
..except for a few that  
double up

# Test: h3ll0 wOrld



# Truth vs Test

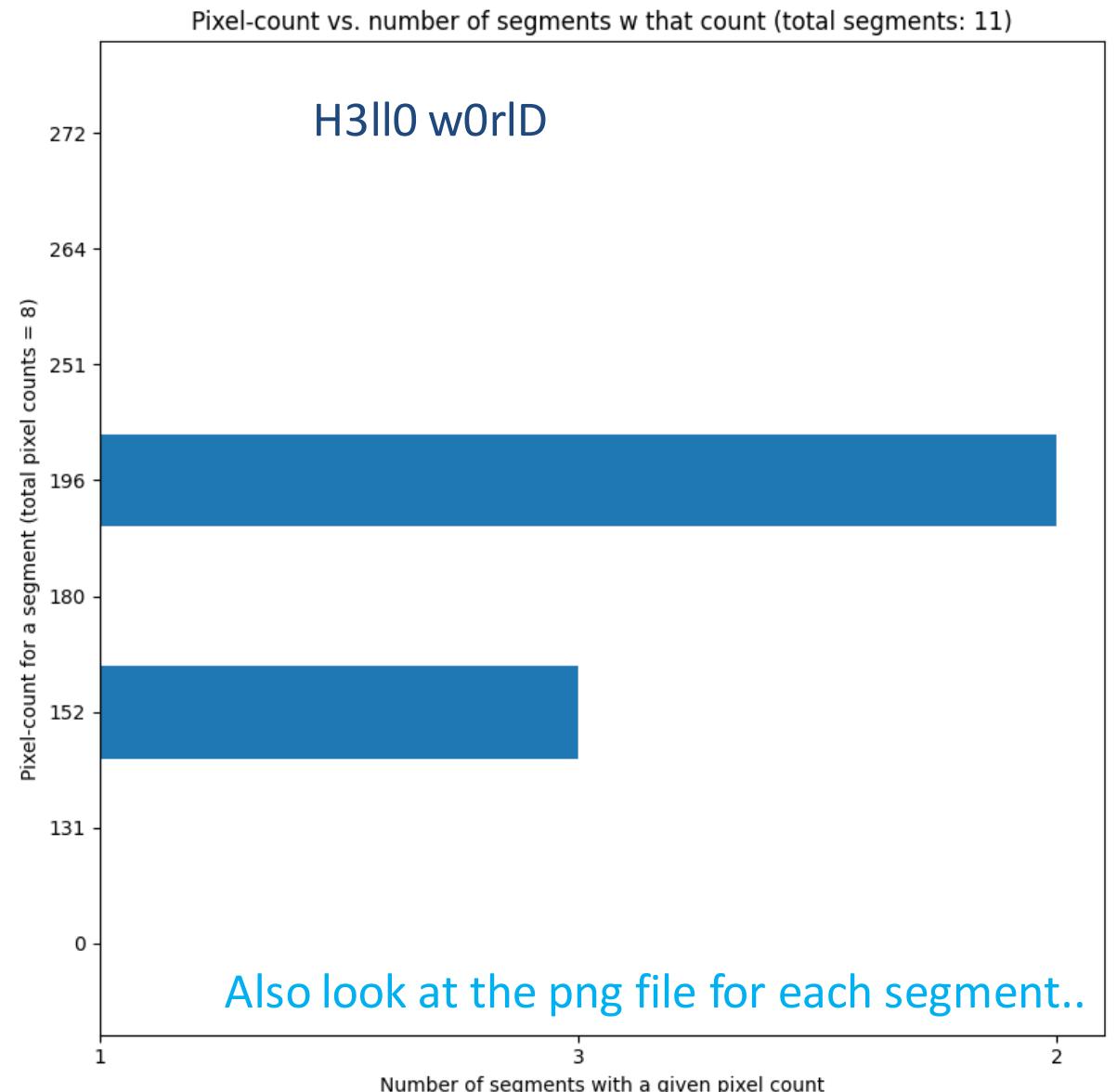
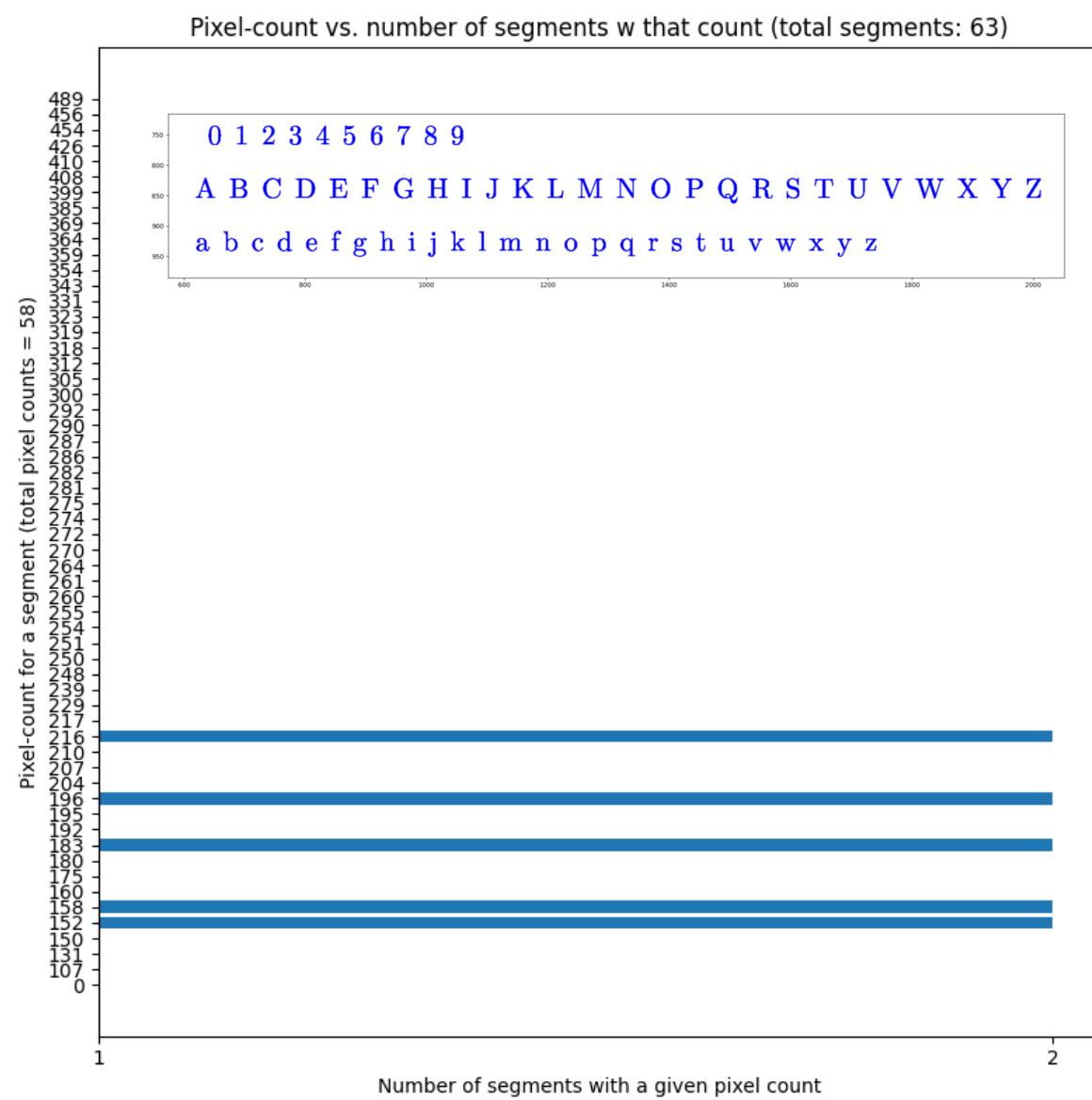


# Truth vs Test



BC Wildfire  
Service

11 = hello world + bg



# Prediction: Earth-mover Distance (EMD)



BC Wildfire  
Service

Traffic plan for moving a pile of stuff... initial and ending formation!

- We used "point" representation!
- Traffic plan = "where to move each piece"
  - Diagram for an "optimal" maneuver!

Method:

1) Compare all possible moves

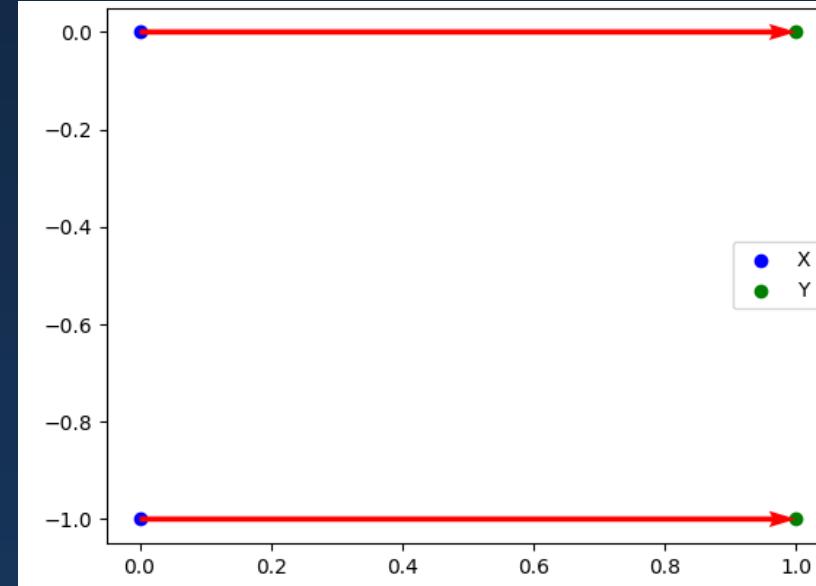
"distance matrix"

2) Choose shortest potential moves first, cross off the items that we add to our plan..

Example:

"move X to Y"

not a lot of combinations...here's one:



```
if __name__ == '__main__':
    a = [[0,1], [0,0]] # [x1, x2], [y1, y2]
    b = [[0,1], [1,1]] # [x1, x2], [y1, y2]
    d, arrows, subdist = dist(a, b)
    dist_plot(a, b, d, arrows, subdist, 0, 1)
```

The other option (X shape) would not be optimal..  
..it is optimal to move the points straight across!  
Dmat: 1, 1,  $\sqrt{2}$ ,  $\sqrt{2}$

[https://en.wikipedia.org/wiki/Earth\\_mover%27s\\_distance](https://en.wikipedia.org/wiki/Earth_mover%27s_distance)

[https://en.wikipedia.org/wiki/Wasserstein\\_metric](https://en.wikipedia.org/wiki/Wasserstein_metric)

# Prediction: Earth-mover Distance (EMD)



BC Wildfire  
Service

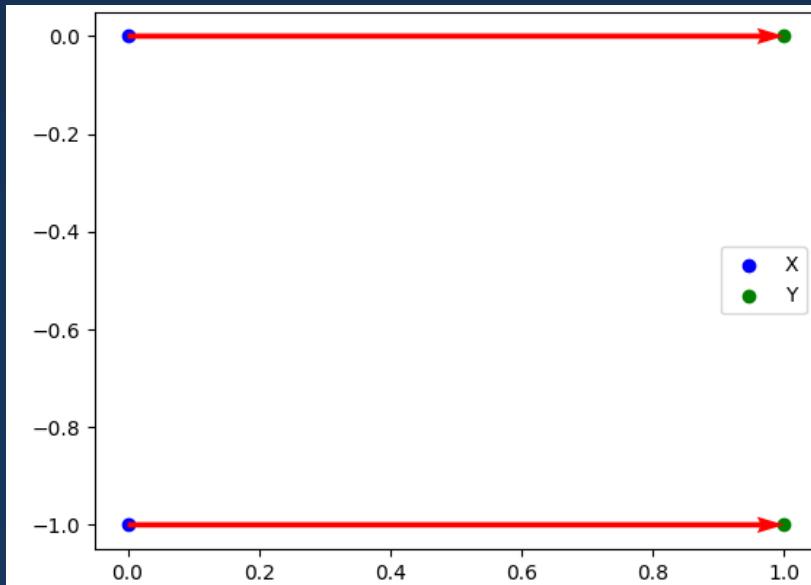
## Method:

1) Compare all possible moves

"distance matrix"

2) Choose shortest potential moves first,  
(for things that aren't moved already)

Add to plan and cross those items off as moved...



Result:  
Compare  
Shapes!

```
def dist(X, Y):
    rho, [x1, y1], [x2, y2] = 0, X, Y # compare two point-sets
    i_n, j_n = len(x1), len(x2)
    subdist, arrows = [], [] # for visualization only

    # element-wise slots for inclusion of each element in traffic plan
    i_f, j_f = [False for i in range(i_n)], [False for i in range(j_n)]

    dm = [[abs(x1[i] - x2[j]) + abs(y1[i] - y2[j]), i, j]
           for i in range(i_n) for j in range(j_n)] # distance matrix
    dm.sort() # sorted distance matrix

    for k in range(len(dm)): # for each dmat element
        d, i, j = dm[k] # dist btwn points, indices of points compared
        if not i_f[i] and not j_f[j]: # if slots open for both elements:
            # add term to distance and close the slots
            i_f[i], i_n, j_f[j], j_n, rho = [True,
                                              i_n - 1,
                                              True,
                                              j_n - 1,
                                              rho + d]

            if d > 0: # record stuff for visualization
                arrows.append([[x1[i], y1[i]], [x2[j], y2[j]]])
                subdist.append(d)

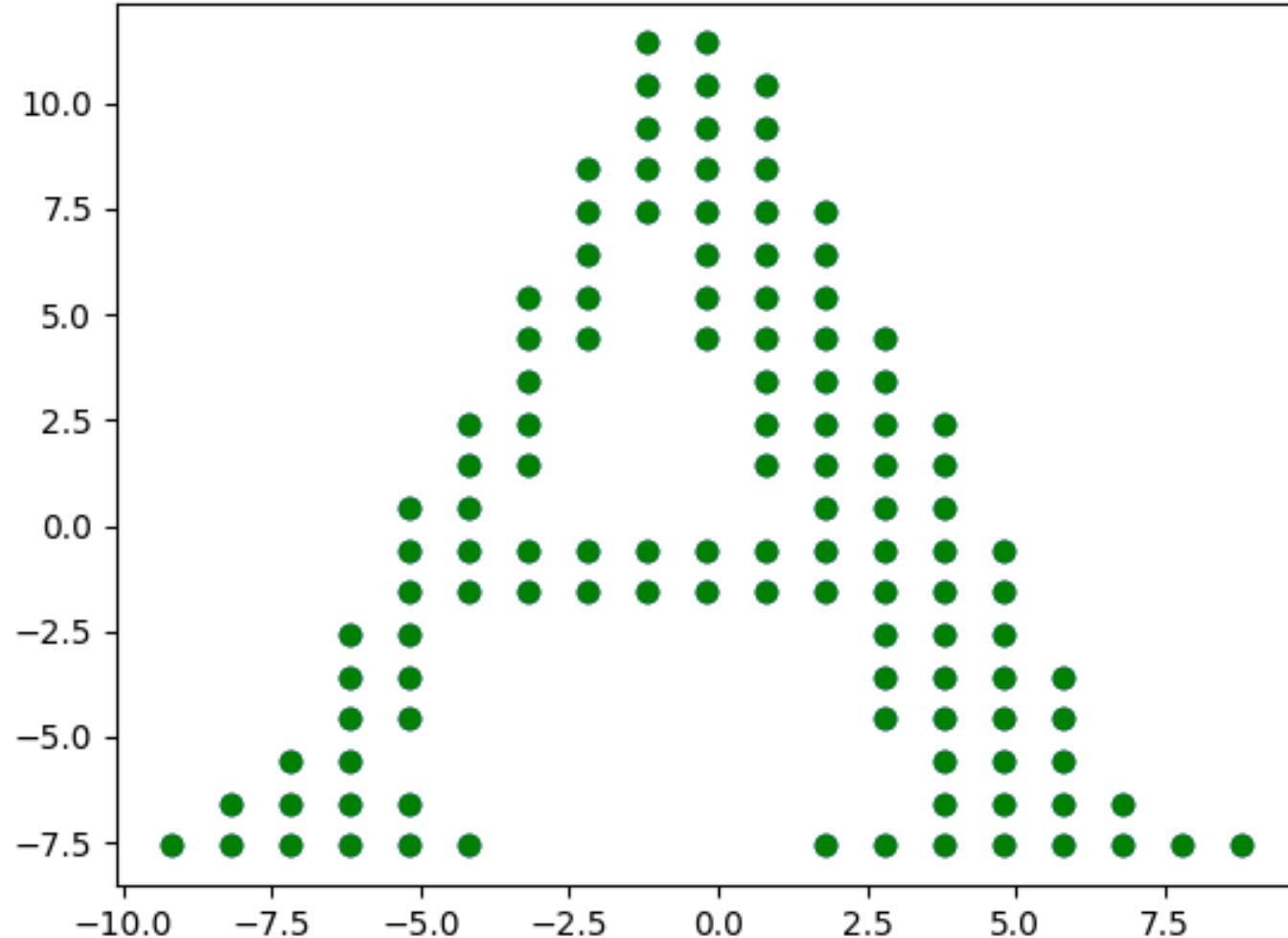
            if i_n * j_n == 0:
                break # ran out of slots for X or Y: stop comparing..

    rho /= min(float(len(x1)), float(len(x2))) # divide by the number of slots
    return rho, arrows, subdist
```

# Prediction: Earth-mover Distance



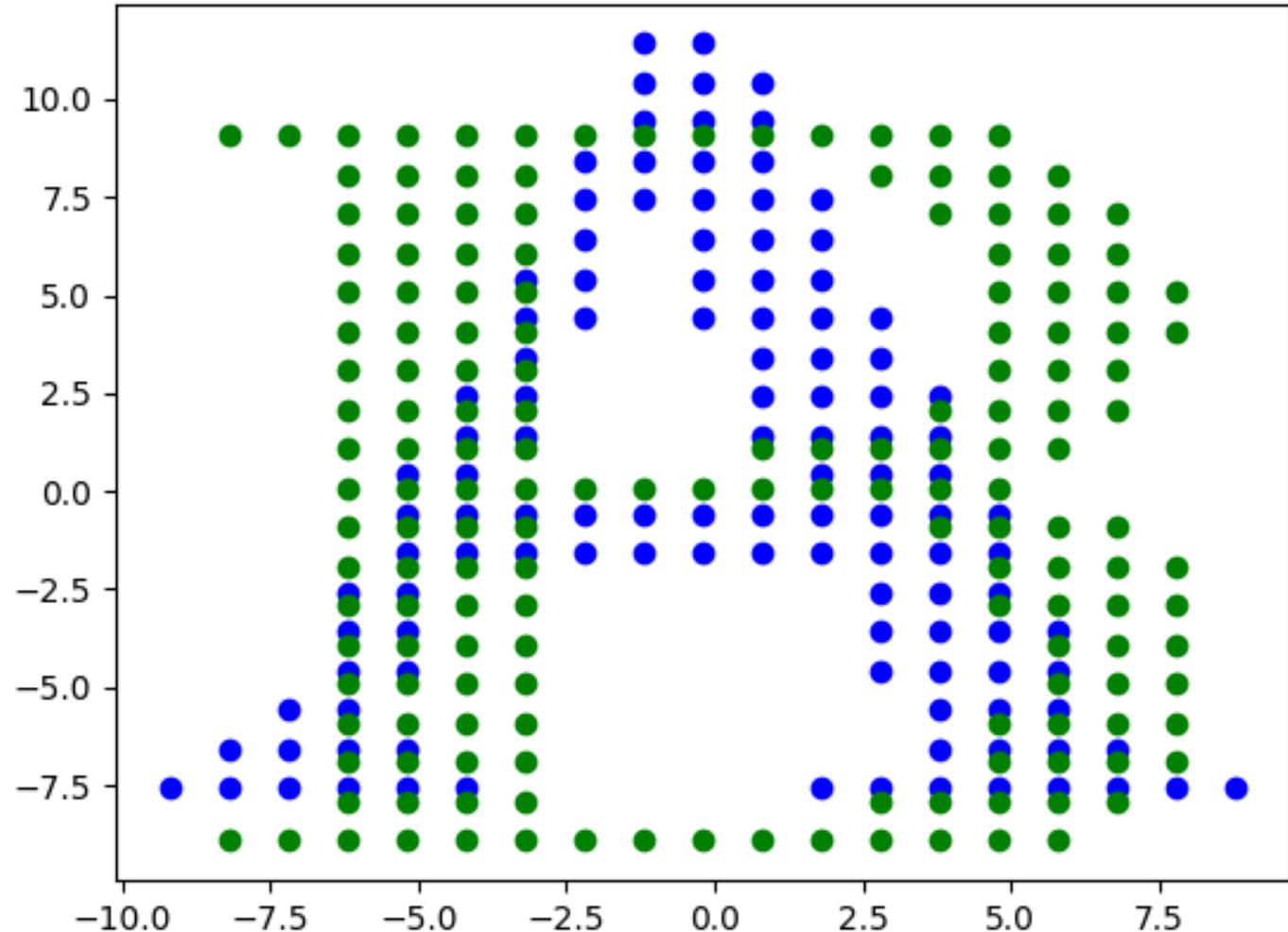
BC Wildfire  
Service



# Prediction: Earth-mover Distance



BC Wildfire  
Service

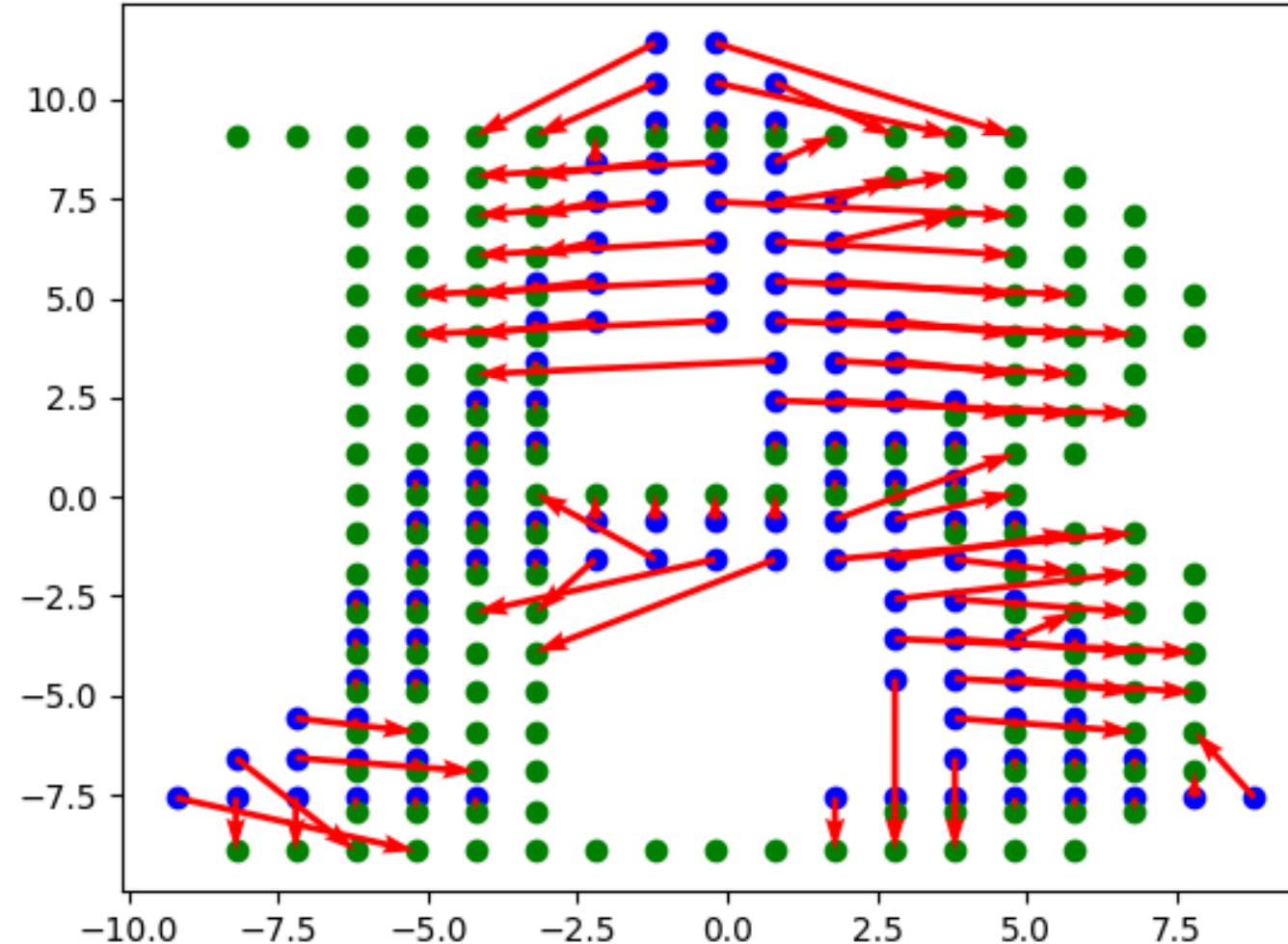


Compare A with B!

# Prediction: Earth-mover Distance



BC Wildfire  
Service

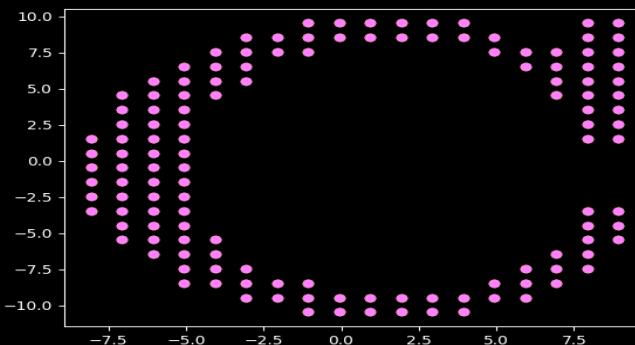
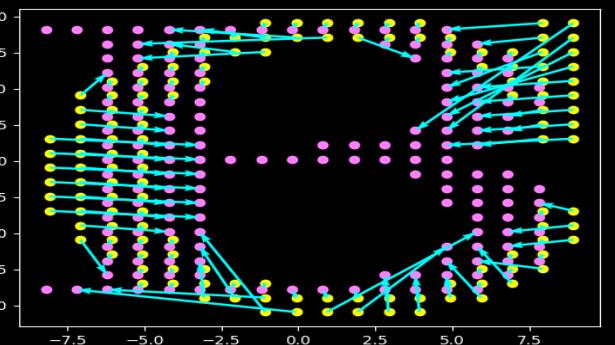
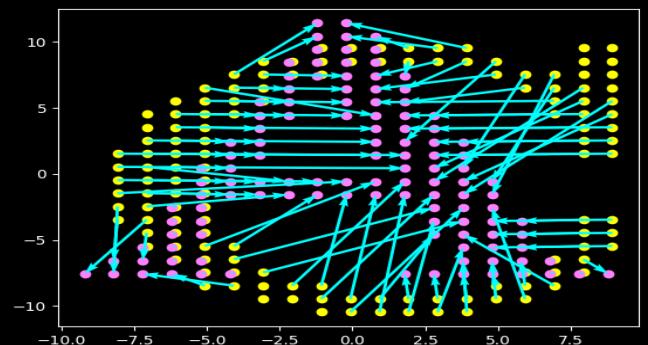
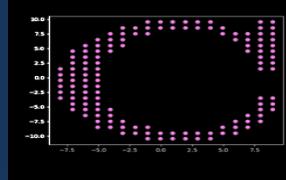
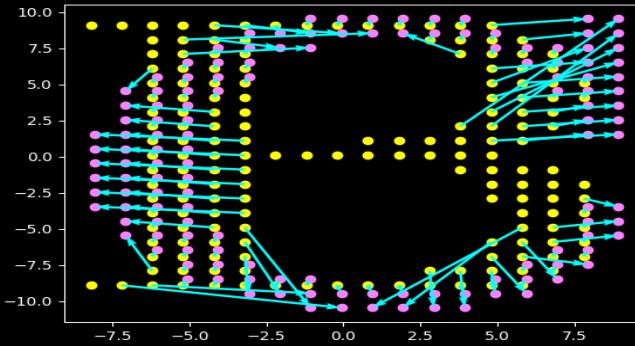
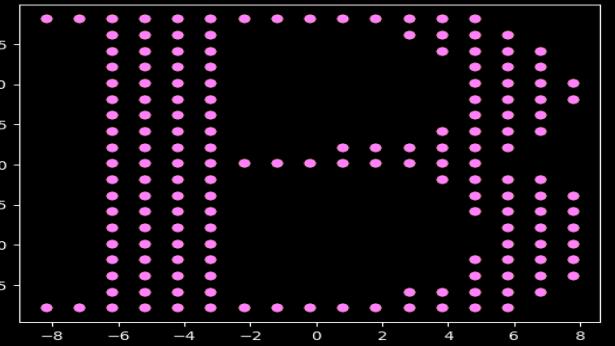
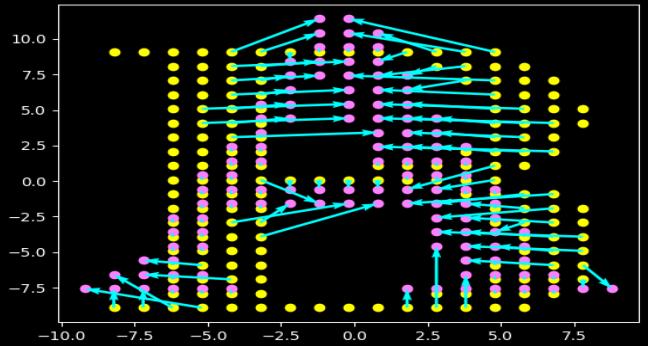
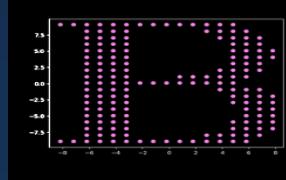
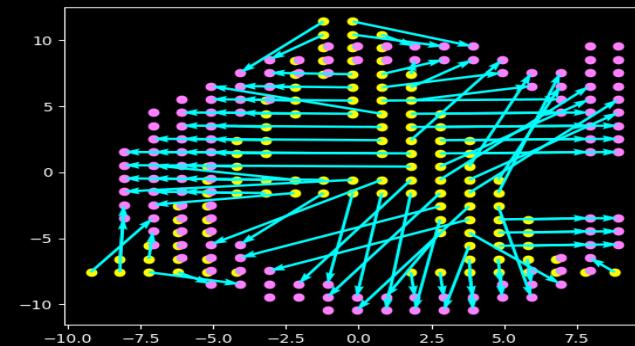
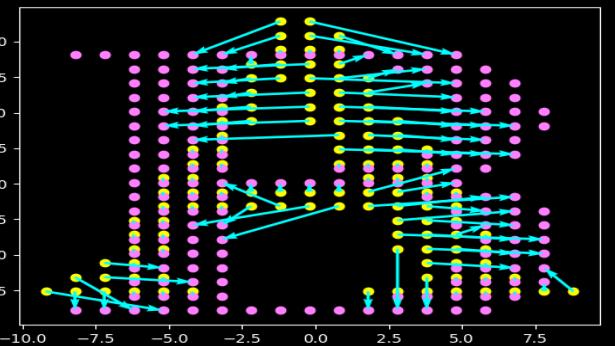
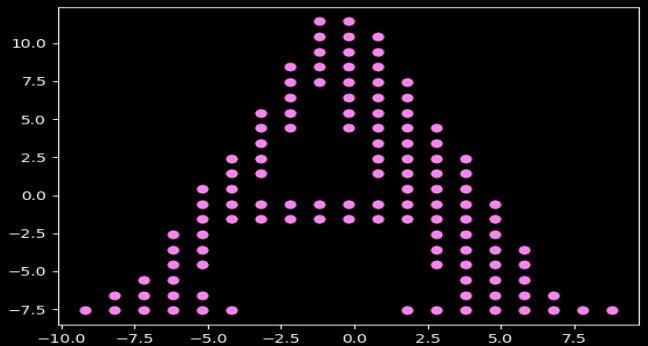
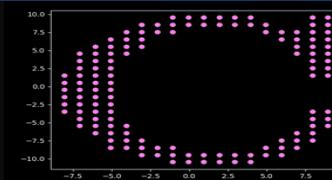
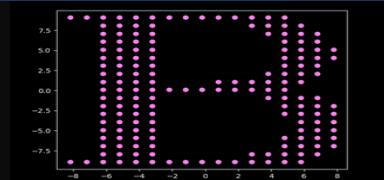
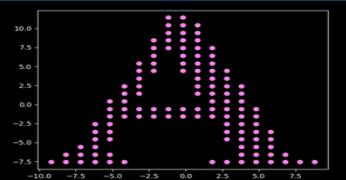


Idea: similar figures take less energy to transform...  
...why? Some tolerance for the figures not being exactly the same



Red:  
"Transport Plan":  
Diagram that says  
what to move  
where

Look at  
relationships  
between letters..  
...matrix!



# How to predict / classify with EMD?

Same as you would with other distance...

Find closest (least distance) "truth" data-point for each test data-point...  
..label the test point accordingly!

"If our test item looks more like an A than any other letters, label it A"

"Pseudocode":

- 1) For each test data (letter):
  - 2) Compare with each truth item (alphabet etc),
  - 3) find the closest match  
(according to EMD "efficient morphing")..
  - 4) Label accordingly... "predict"

Run in parallel!!!!!! Because...

```
def predict_i(pi): # for pi in range(len(test_points)):  
    print("predict_i", pi)  
    # print(test_files[pi])  
    p = normalize(test_points[pi]) # didn't normalize before to get absolute centroid  
  
    min_d, min_i = sys.float_info.max, None # closest truth value  
  
    for i in range(len(truth_points)):  
        t = normalize(truth_points[i])  
        d, arrows, subdist = dist(p, t)  
        print("rho", d, "i", i)  
  
        if pi == 0:  
            dist_plot(p, t, d, arrows, subdist, pi, i)  
  
        if d < min_d: # found a better match  
            min_d, min_i = d, i  
  
        if min_d == 0.: # perfect match, stop looking..  
            break  
  
    print("min_i", min_i)  
    prediction = truth_labels[min_i]  
    result = [test_centroids[pi], prediction]  
    print(result)  
    return(result)  
    print("point", pi, "of", len(test_points))  
  
  
print("truth_points", truth_points)  
use_parfor = True  
if use_parfor:  
    predictions = parfor(predict_i, range(len(test_points)))  
else:  
    predictions = [predict_i(pi) for pi in range(len(test_points))]  
  
print("prediction", predictions)  
  
plt.figure()  
for p in predictions:  
    plt.plot(p[0][1], -p[0][0])  
    plt.text(p[0][1], -p[0][0], p[1])  
plt.show()  
plt.savefig("prediction.png")
```

# Prediction for h3ll0 w0rlD..



truth ['h', '3', 'l', 'l', '0', 'w', '0', 'r', 'l', 'D']

prediction ['3', 'D', 'l', '0', 'h', 'r', '0', 'l', 'l', 'w']

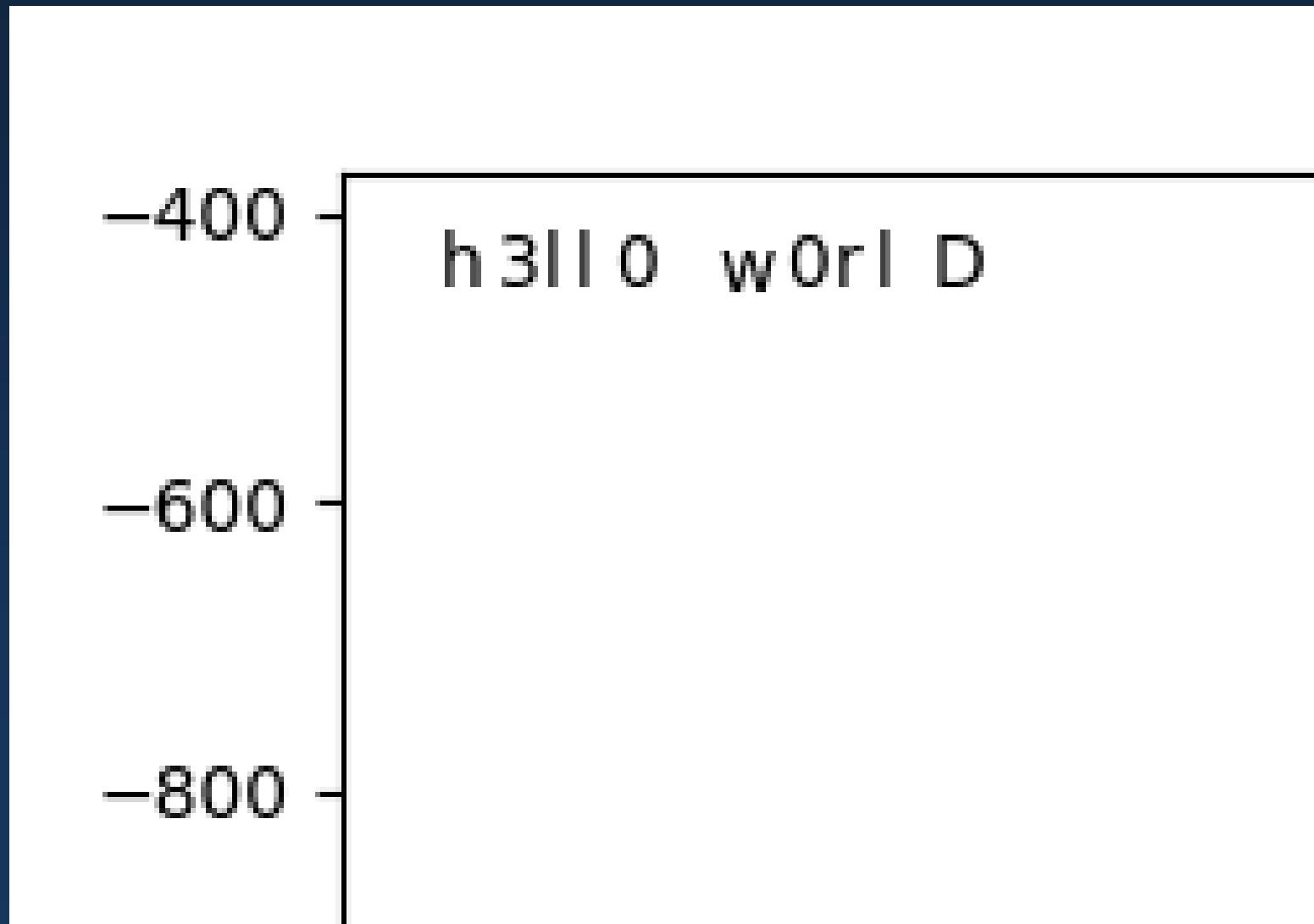
The letters aren't the same height...  
....scanning right-left top-bottom doesn't pull them in order.

- For truth data set scan lines (in the segmentation code...)
- For predictions, plot the predicted value at the centroid location

# Predictions: plot at centroid locations..



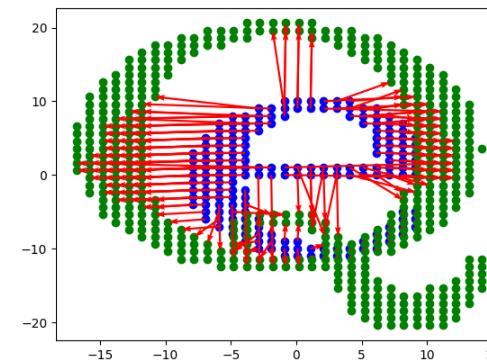
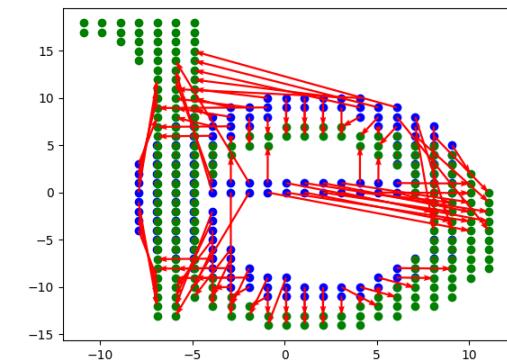
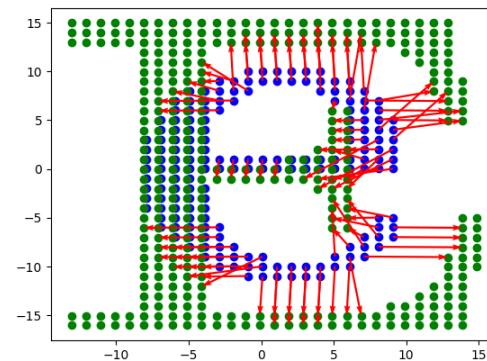
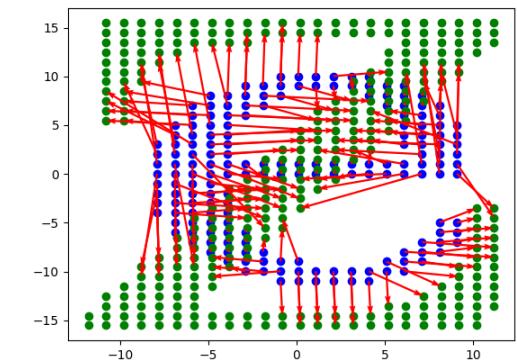
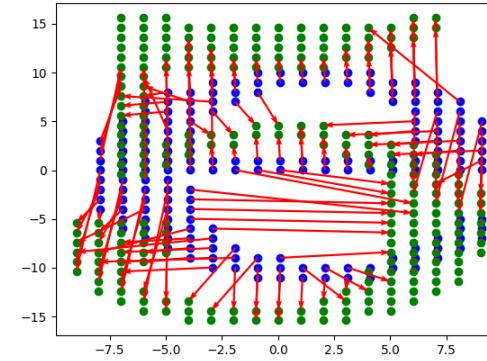
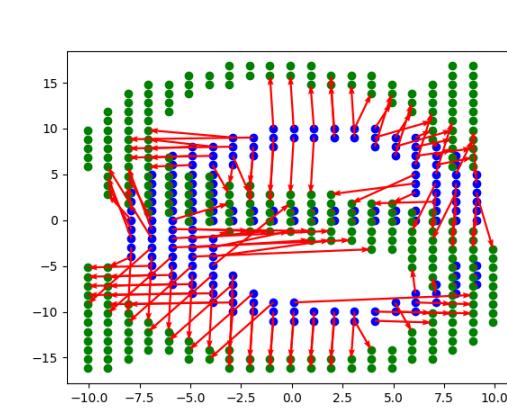
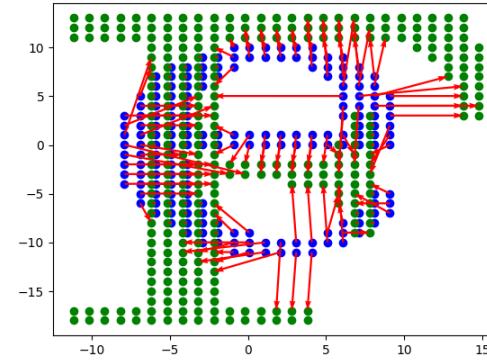
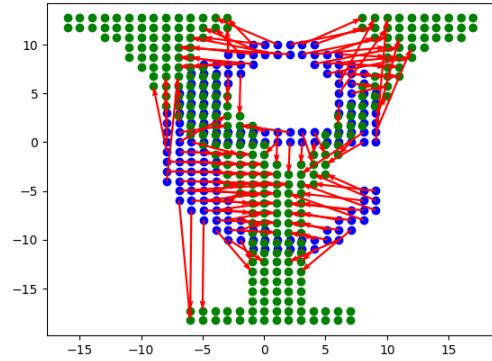
BC Wildfire  
Service



# Prediction: example... Match onto "e" ..



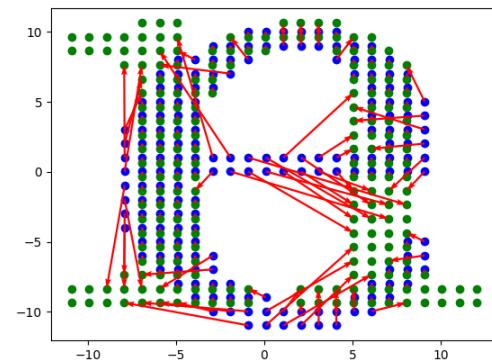
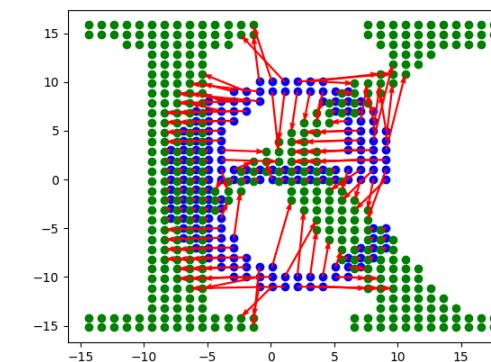
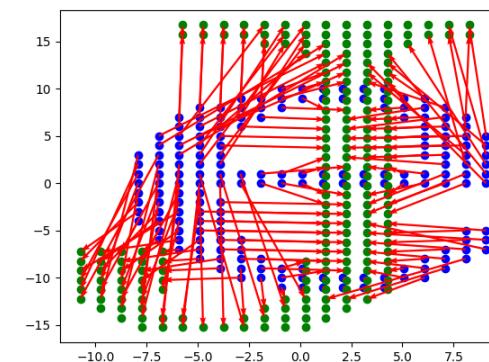
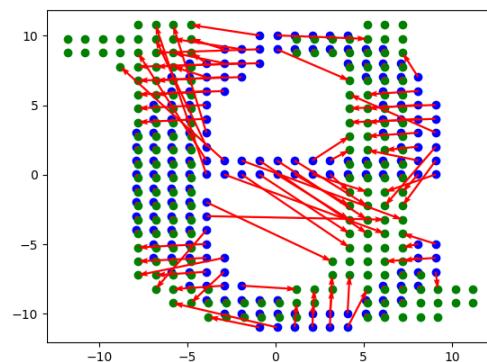
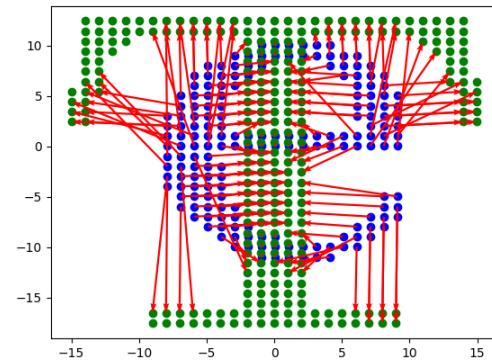
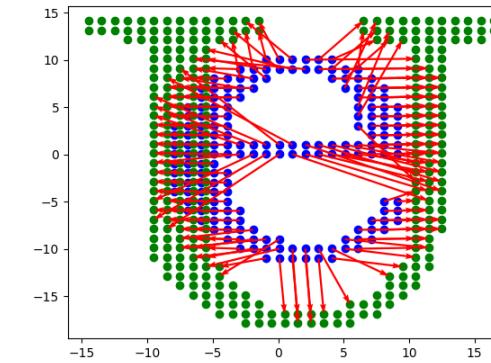
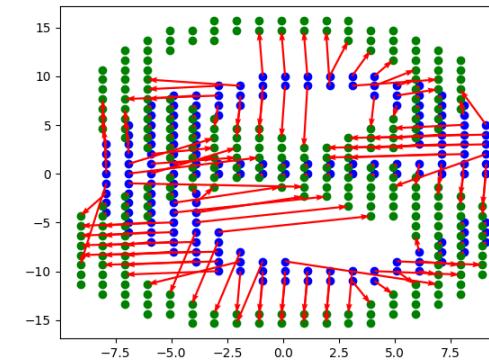
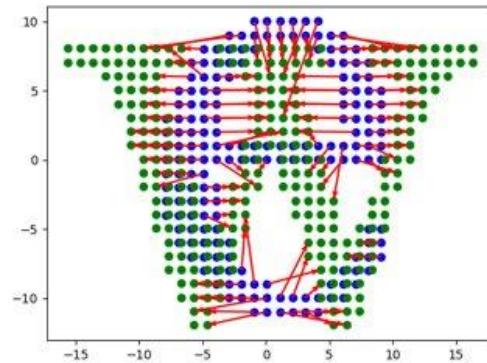
BC Wildfire  
Service



# Prediction: example... Match onto "e"...



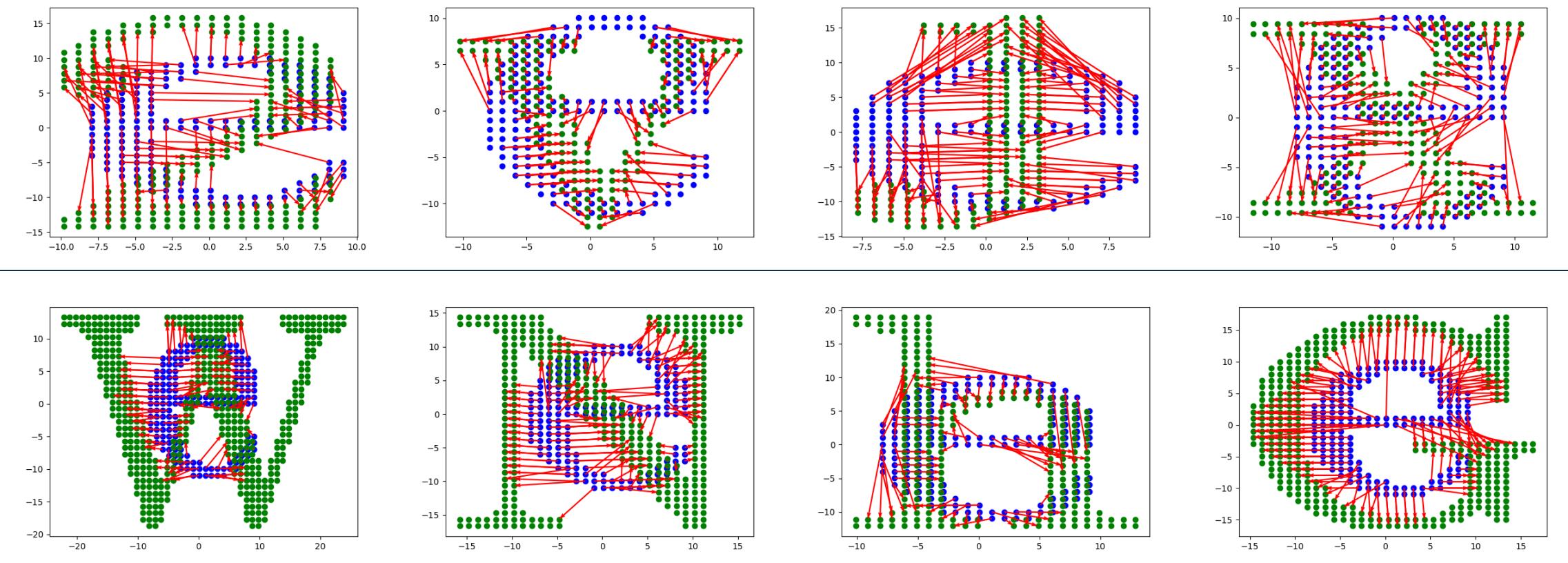
BC Wildfire  
Service



# Prediction: example... Match onto "e"...



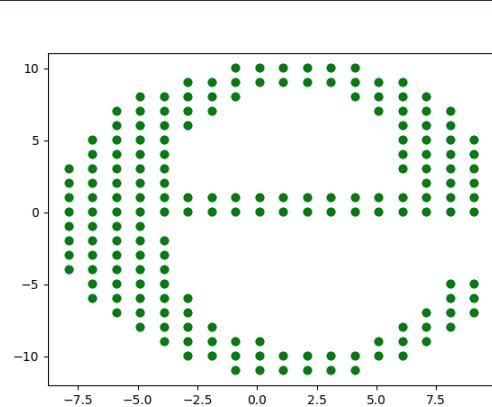
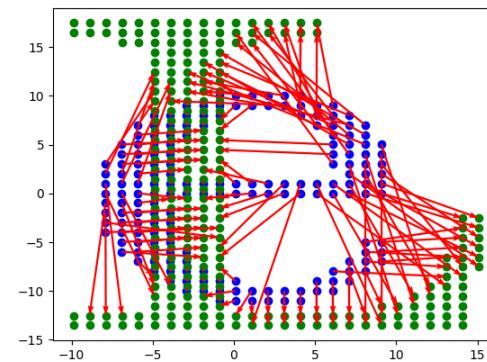
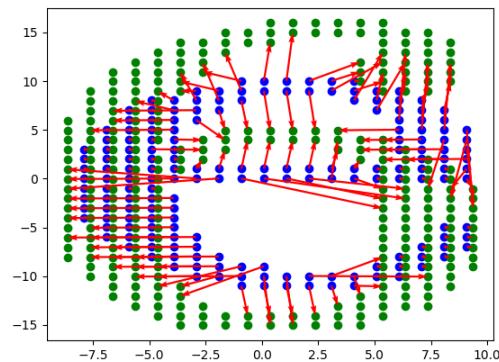
BC Wildfire  
Service



# Prediction: example... Match onto "e" ..



BC Wildfire  
Service



Perfect overlap..

..Got a good match, can  
stop now!

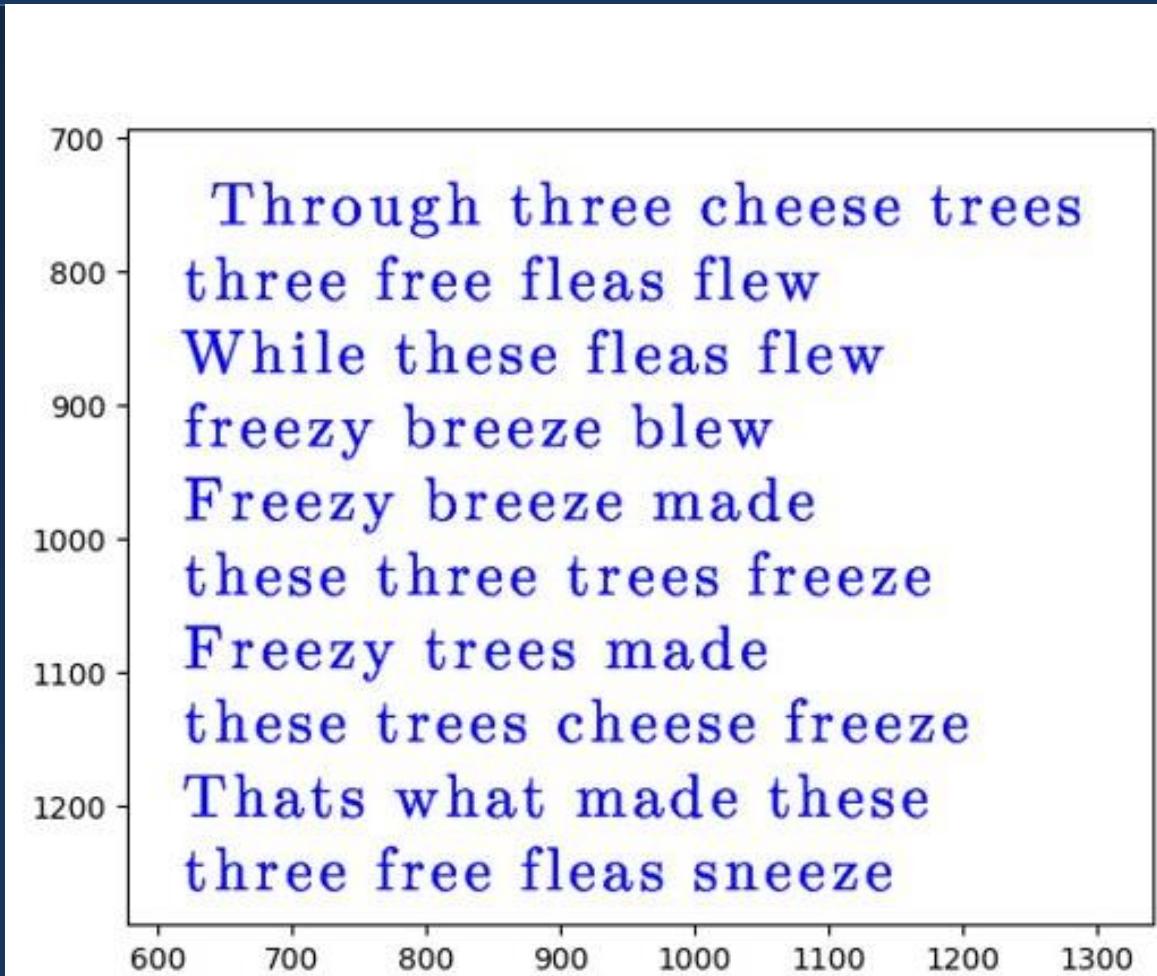
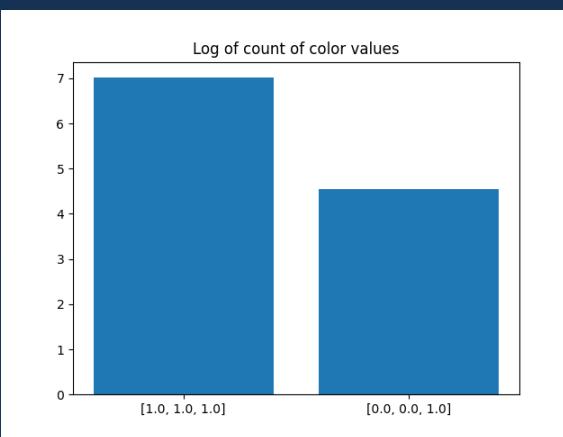
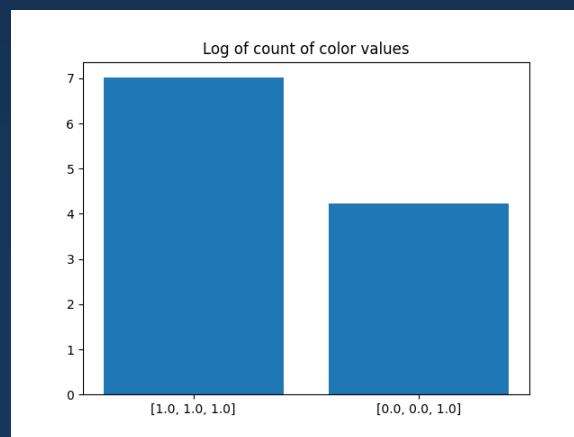
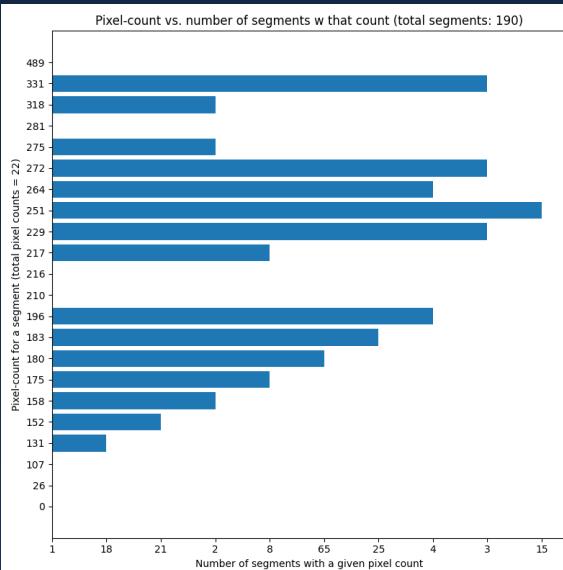
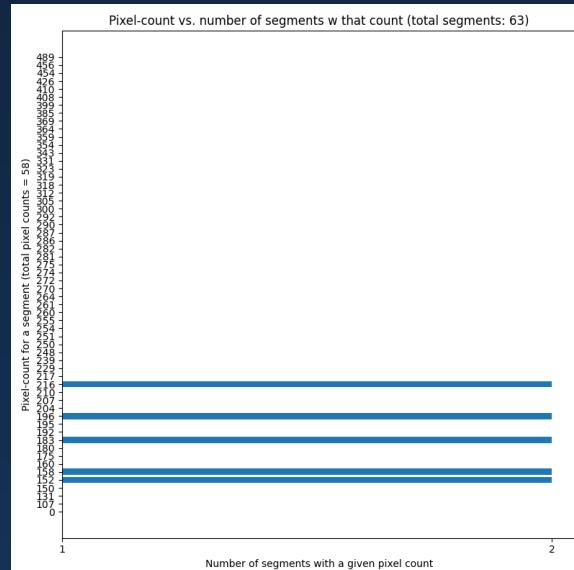
Good news:

This will still work if the "e" we were matching was a slightly different shape from the "e" in the "truth" set....

# Scaling up..

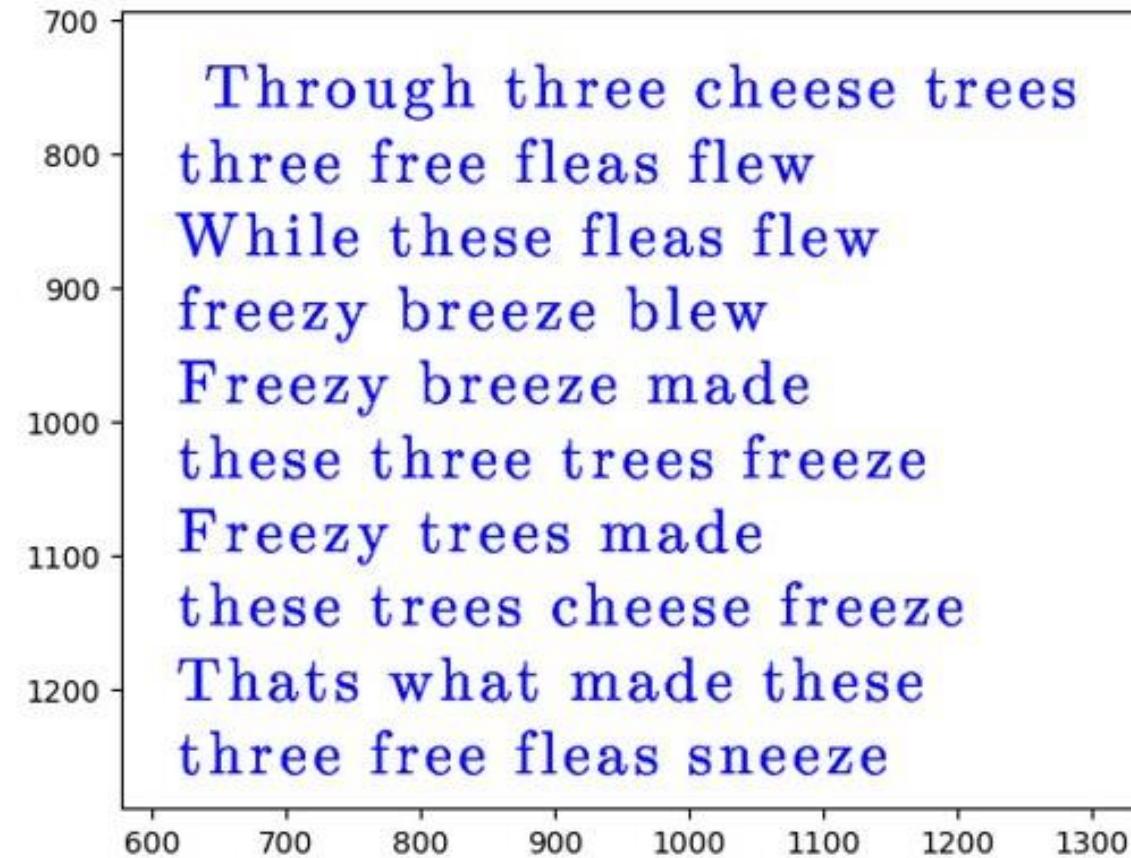


BC Wildfire  
Service

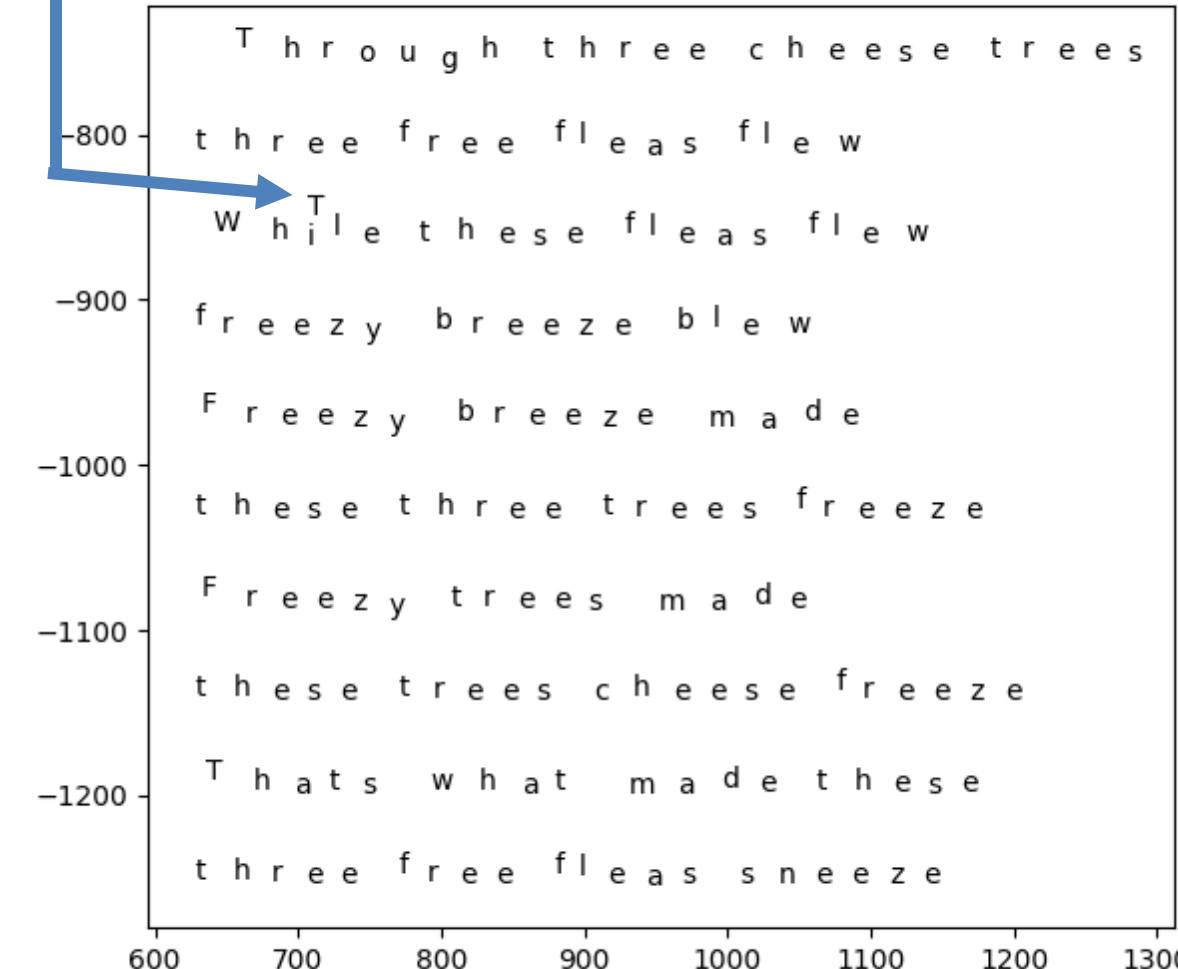


# Prediction

What's that?



Need a segment size threshold..



# Earth Mover + Predictive Models.. Applications



- Compare patterns / shapes!
  - By morphing them
  - Makes sense for many kinds of data..
- Image classification.. Sensor fusion for environment monitoring
  - Land
  - Atmosphere
  - LiDAR, Multispectral, Radar, Hyperspectral etc...(or combos)..
- Resource allocation / distribution..
  - Cost efficiency for maneuvers (planes, helicopters etc.)...
    - route finding for groups, rather than just individuals!

# PyTesseract example!!!



Very short script for OCR using OpenCV and Tesseract (thanks Ken Fullerton for the share!)

```
# text recognition

import cv2
import pytesseract
# read image
im = cv2.imread('./test3.jpg')
# configurations
config = ('-l eng --oem 1 --psm 3')
# pytesseract
text = pytesseract.image_to_string(im, config=config)

# print text
text = text.split('\n')
text
```



Questions?

[Ashlin.Richardson@gov.bc.ca](mailto:Ashlin.Richardson@gov.bc.ca)



**BC Wildfire  
Service**