

Image caption generator application

A PROJECT REPORT

Submitted by

Vineet Sharma (19MIM10001)

Thevapraakash P (19MIM10003)

Shivanshu Rajput (19MIM10060)

Shrey Asthana (19MIM10066)

Kesavan R (19MIM10086)

BONAFIDE CERTIFICATE

Certified that this project report titled "**Image caption generator application** " is the Bonafide work of "**VINEET SHARMA (19MIM10001), THEVAPRAKASH P (19MIM10003), SHIVANSHU RAJPUT (19MIM10060), SHREY ASTHANA (19MIM10066)** and **KESAVAN R (19MIM10086)**" who carried out the project work under my supervision, certified further that to the best of my knowledge the work reported here does not form part of any other project / research work on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

ACKNOWLEDGEMENT

Primarily, I would like to thank the Lord Almighty for his presence and immense blessings throughout the project work.

I wish to express my heartfelt gratitude to our instructor Pradeepthi Duggaraju, for continually guiding and actively participating in our project, giving valuable suggestions to complete the project work.

I would like to thank all the technical and teaching staff of the smart internz and smart bridge, who extended directly or indirectly all support.

Last but not the least, I am deeply indebted to my parents who have been the greatest support while I worked day and night for the project to make it a success.

Index

Serial number	Topic
1	Introduction
2	Problem statement
3	Solution to the problem
4	Literature survey
5	Experimental analysis
6	Conclusion
7	Application
8	Future scope
9	Code snippet
10	Bibliography

INTRODUCTION

The world around us is incredibly beautiful, full of wonders and to see them we have been provided with the eyes. And it is really disappointing to see many people with partially and complete blindness. It is our effort to improvise their life by providing the description of the image.

It is a tedious task to describe the content of the image captured using the properly English statement but it could have a significant impact by helping visually impaired people better understand their surroundings.

In the modern era of technology most of the cell phones have the built-in cameras to capture the surrounding, making it convenient to visualise the surrounding for the visual impaired person so that they can be benefitted and helped to overcome the problem that they have faced frequently.

PROBLEM STATEMENT

The problem is that the blind person must depend on someone or something to explore the

surrounding. The need of stick is mandatory for a blind person. We try to reduce the burden of

the blind person by providing him the description of the surrounding by just seeing the

surrounding in the form of an image.

SOLUTION TO THE PROBLEM

We are creating a web application where the user selects the image and the image is fed into the

model that is trained and generated caption will be displayed on the webpage. This generated

caption could be read aloud in the future for better aid to the visually impaired people.

LITERATURE SURVEY

Existing problem

For a blind person, it is really disturbing to not be able to see the surrounding. In fact, they have to

depend on the other senses for perceiving the environment. The only way to overcome the

blindness is the transplant of the retina by which they can be able to see the neighboring thing but

that is not affordable to each and everyone.

Another problem is that if the person is suffering from squint, then the transplant would also

does not affect the situation of the person.

Proposed solution

A method to somehow visualize the surrounding to the visual impaired person is our priority in

this project. We are creating a web application where the user selects the image and the image

is fed into the model that is trained and generated caption will be displayed on the webpage.

This generated caption could be read aloud in the future for better aid to perceive the

environment for the visually impaired people.

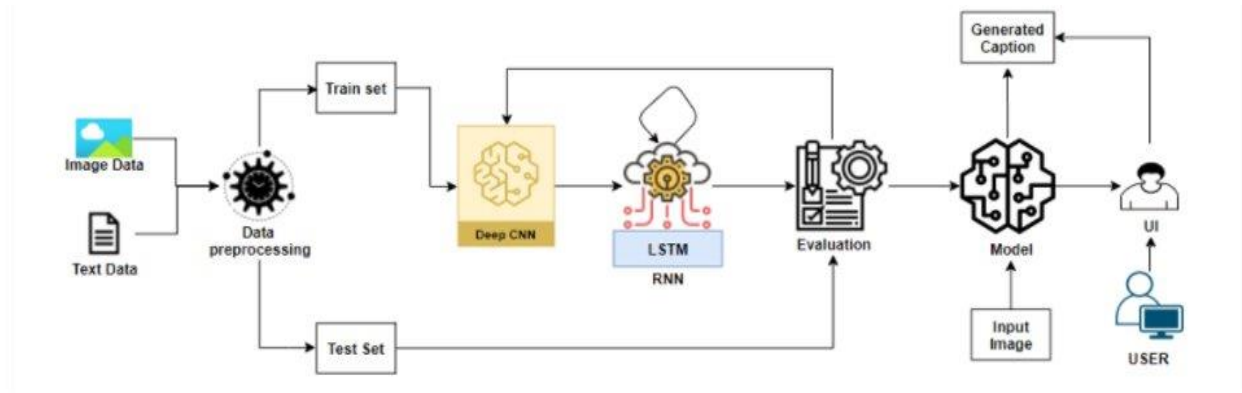
EXPERIMENTAL ANALYSIS

To accomplish our project whose main objective is to describe the image in the best form to the viewer using the CNN-RNN model for the visual impaired person. We need

complete the task enlisted below: -

- Data Collection
 - Collect the dataset or create the dataset
- Data Preprocessing.
 - Import required Libraries
 - Extract features from each photo in the directory
 - Processing the text data or descriptions
- Model Building
 - Import the model building Libraries
 - Loading dataset for training the model
 - Tokenizing the Vocabulary
 - Define the Model
 - Define the CNN-RNN Model
 - Configure the Learning Process
 - Training the model
 - Save the Model
 - Testing the Model
- Application Building
 - Create an HTML file
 - Build Python Code

System architecture for the project



Hardware requirement

The laptop with a good broadband is required to accomplish the project.

Laptop must have a bare minimum of 8 gpu.

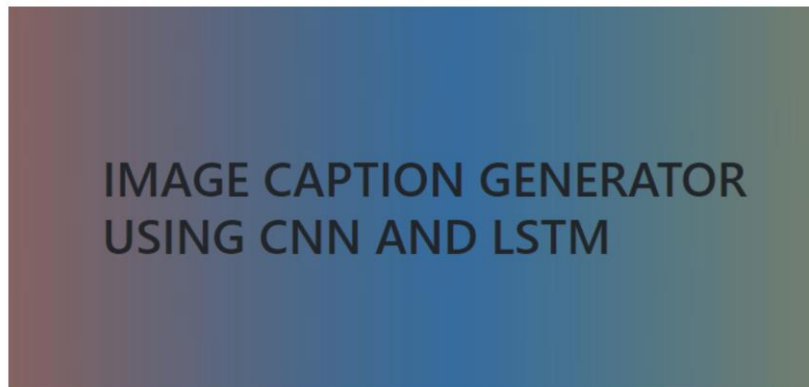
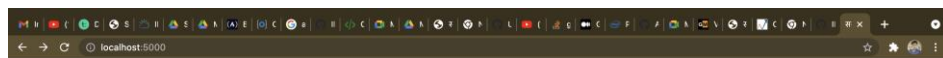
Software requirement

- Python 3.9
- Anaconda navigator
- TensorFlow version 1.14.0
- Keras 2.2.4
- Flask
- And other python libraries like NumPy, pandas, OpenCV, matplotlib and many more.

CONCLUSION

The following are the images of our project in which we are producing the captions for the

images uploaded.



THE PROBLEM STATMENT



THE PROBLEM STATEMENT

For a machine to be able to automatically describe objects in an image along with their relationships or the actions being performed using a learnt language model is a challenging task, but with massive impact in many areas. Being able to automatically describe the content of an image using properly formed English sentences is a challenging task, but it could have great impact by helping visually impaired people better understand their surroundings. Most modern mobile phones are able to capture photographs, making it possible for the visually impaired to make images of their environments. These images can then be used to generate captions that can be read out loud to the visually impaired, so that they can get a better sense of what is happening around them.

WHAT IS CNN

Swag shoiuldgoitch literally meditation subway tile
tumblr cold-pressed. Gastropub street art beard
dreamcatcher neutra, ethical XOXO lumbersexual.

WHAT IS LSTM

Swag shoiuldgoitch literally meditation subway tile
tumblr cold-pressed. Gastropub street art beard
dreamcatcher neutra, ethical XOXO lumbersexual.

WHAT IS IMAGE CAPTION
GENERATOR

Swag shoiuldgoitch literally meditation subway tile
tumblr cold-pressed. Gastropub street art beard
dreamcatcher neutra, ethical XOXO lumbersexual.



NAME

For a machine to be able to automatically describe objects in an image along with their relationships or the actions being performed using a learnt language model is a challenging task, but with massive impact in many areas. Being able to automatically describe the content of an image using

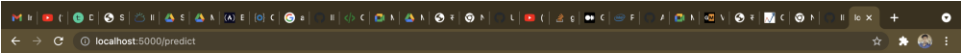


NAME

For a machine to be able to automatically describe objects in an image along with their relationships or the actions being performed using a learnt language model is a challenging task, but with massive impact in many areas. Being able to automatically describe the content of an image using properly formed English sentences is a challenging task, but it could have great impact by helping visually impaired people better understand their surroundings. Most modern mobile phones are able to capture photographs, making it possible for the visually impaired to make images of their environments. These images can then be used to generate captions that can be read out loud to the visually impaired, so that they can get a better sense of what is happening around them.

Choose Upload 

File : 99679241_adc853a5c0.jpg



startseq pelican

APPLICATION

The following application could be used for better understanding of the uploaded images. With the voice assistance we can make a better application for the blind person to visualize the surrounding in a much better way.

FUTURE SCOPE

There is always a scope of improvement in each and everything created. On the same principle

there might be the future improvement in this project also. The usage of voice could be our next

step to enhance the project. By the addition of the vocals would be boon to the blind person as

they can hear the description of the surroundings.

Next, we could launch the versions of this web application onto the different operating systems

like android and iOS.

CODE SNIPPET

```
1 from pickle import load
2 from numpy import argmax
3 from tensorflow.keras.preprocessing.sequence import pad_sequences
4 from tensorflow.keras.applications.vgg16 import VGG16
5 from tensorflow.keras.preprocessing.image import load_img
6 from tensorflow.keras.preprocessing.image import img_to_array
7 from tensorflow.keras.applications.vgg16 import preprocess_input
8 from tensorflow.keras.models import Model
9 from tensorflow.keras.models import load_model
10 import os
11 from flask import Flask, render_template, request
12 from werkzeug.utils import secure_filename
13 from event.pywsgi import WSGIServer
14
15 app = Flask(__name__)
16 @app.route('/')
17 def home():
18     return render_template("index.html")
19
20 @app.route('/predict', methods = ['GET', 'POST'])
21 def upload():
22     if request.method == "POST":
23         f = request.files["image"]
24         print('current path', basepath)
25         basepath = os.path.dirname(__file__)
26         print('current path', basepath)
27         filepath = os.path.join(basepath, "uploads", f.filename)
28         print('upload folder is', filepath)
29         f.save(filepath)
30         text = model.predict(filepath)
31         return text
32
33 def extract_features(filename):
34     print('Features extracted')
35     model = VGG16()
36     model.layers.pop()
37     model = Model(inputs = model.inputs, outputs = model.layers[-1].output)
38     image = load_img(filename, target_size=(224, 224))
39     print('image loaded')
40     image = img_to_array(image)
41     image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
42     image = preprocess_input(image)
43     feature = model.predict(image, verbose=0)
44     print('model predicted')
45     return feature
46
47 def word_for_id(integer, tokenizer):
48     for word, index in tokenizer.word_index.items():
49         if index == integer:
50             return word
51     return None
52
53 def generate_desc(model, tokenizer, photo, max_length):
54     print('generate description')
55     in_text = "startseq"
56     for i in range(max_length):
57         sequence = tokenizer.texts_to_sequences([in_text])[0]
58         sequence = pad_sequences([sequence], maxlen=max_length)
59         print('sequence')
60         yhat = model.predict([photo, sequence], verbose=0)
61         yhat = argmax(yhat)
62         word = word_for_id(yhat, tokenizer)
63         if word is None:
64             break
65         in_text += ' ' + word
66         if word == 'endseq':
67             break
68     print(in_text)
69     return in_text
70
71 def modelpredict(filepath):
72     tokenizer = load_model("../Users/shrey/SmartBridge/tokenizer.pkl", 'rb')
73     max_length = 34
74     model = load_model("../Users/shrey/SmartBridge/caption.h5")
75     print('model loaded')
76     photo = extract_features(filepath)
77     description = generate_desc(model, tokenizer, photo, max_length)
78     return description
79
80 if __name__ == "__main__":
81     app.run(debug = True)
```

Usage

Here you can get help of any object by pressing Cmd+I in front of it, either on the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in **Preferences > Help**.

New to Spyder? Read our [tutorial](#)

Variable explorer Help Plots Files

Console 1/A

To enable them in other operations, rebuild TensorFlow with the appropriate compiler flags.

2021-07-30 18:54:20.157366: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:176] None of the MLIR Optimization Passes are enable (registered 2)

127.0.0.1 -- [30/Jul/2021 18:54:21] "POST /predict HTTP/1.1" 200 -

127.0.0.1 -- [30/Jul/2021 18:54:46] "POST /predict HTTP/1.1" 200 -

WARNING:tensorflow: Out of the last 5 calls to <function Model.make_predict_function.<locals>.predict_function at 0x7fa579485280> triggered tf.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop, (2) passing tensors with different shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.function has experimental_relax_shapes=True option that relaxes argument shapes that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/functioncontrolling_retracing and https://www.tensorflow.org/api_docs/python/tf/function for more details.

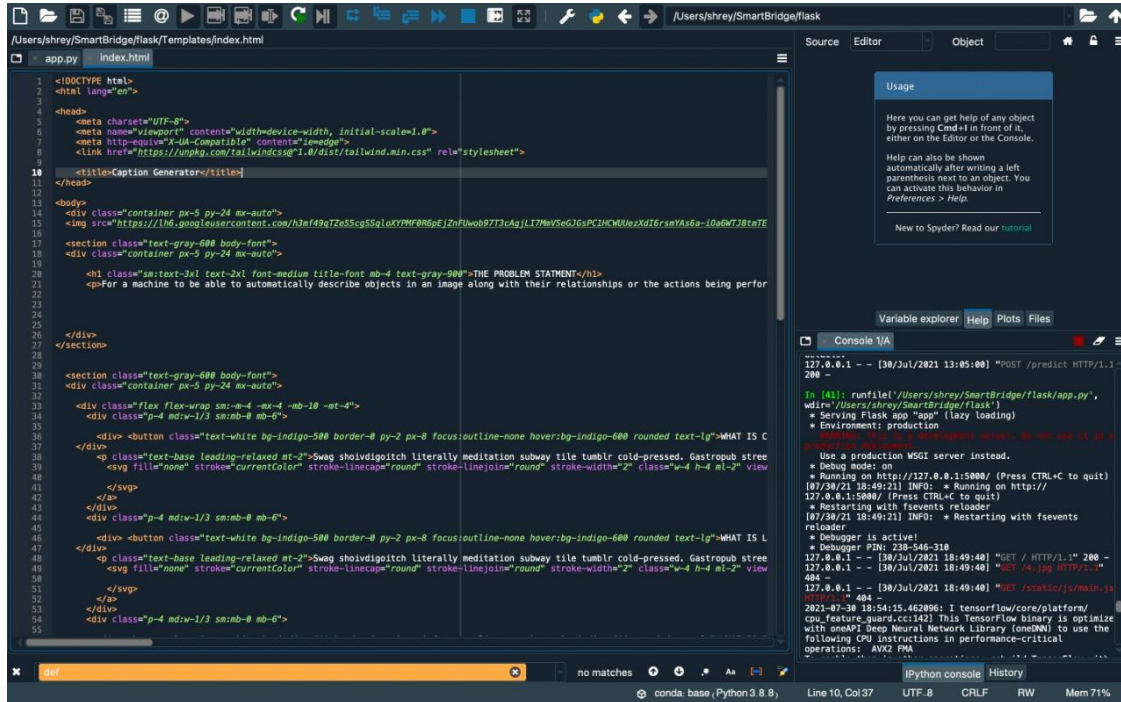
WARNING:tensorflow: Out of the last 6 calls to <function Model.make_predict_function.<locals>.predict_function at 0x7fa5f422df70> triggered tf.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop, (2) passing tensors with different shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.function has experimental_relax_shapes=True option that relaxes argument shapes that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/functioncontrolling_retracing and https://www.tensorflow.org/api_docs/python/tf/function for more details.

127.0.0.1 -- [30/Jul/2021 18:59:17] "POST /predict HTTP/1.1" 200 -

Python console History

LSP Python: ready conda base (Python 3.8.8) Line 23, Col 33 ASCII LF RW Mem 71%

This is app.py file.



The screenshot displays the Spyder IDE interface. The main editor window shows the `app.py` file, which contains HTML code for a web application. The code includes a meta tag for charset, viewport, and a link to a stylesheet. It also features a section for a "Caption Generator" and a form with a text input and a submit button. The console window on the right shows the output of the application, including the Flask app's startup message and the results of a POST request to the `/predict` endpoint.

```
1 <!DOCTYPE html>
2 <html lang="en">
3
4 <head>
5   <meta charset="UTF-8">
6   <meta name="viewport" content="width=device-width, initial-scale=1.0">
7   <meta http-equiv="X-UA-Compatible" content="ie=edge">
8   <link href="https://unpkg.com/tailwindcss@1.0/dist/tailwind.min.css" rel="stylesheet">
9
10  <title>Caption Generator</title>
11 </head>
12
13 <body>
14   <div class="container px-5 py-24 mx-auto">
15     
16   </div>
17   <div class="text-gray-600 body-font">
18     <div class="container px-5 py-24 mx-auto">
19
20       <h1 class="sm:text-3xl text-2xl font-medium title-font mb-4 text-gray-900">THE PROBLEM STATEMENT</h1>
21       <p>For a machine to be able to automatically describe objects in an image along with their relationships or the actions being performed</p>
22
23     </div>
24   </div>
25
26   <div class="text-gray-600 body-font">
27     <div class="container px-5 py-24 mx-auto">
28
29       <div class="flex flex-wrap sm:-mx-4 -mx-4 -mb-10 -mt-6">
30         <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
31
32           <div>
33             <button class="text-white bg-indigo-500 border-0 py-2 px-8 focus:outline-none hover:bg-indigo-600 rounded text-lg">WHAT IS C</button>
34           </div>
35           <p>class="text-base leading-relaxed mt-2">Swag shoidvigolitch literally meditation subway tile tumblr cold-pressed. Gastropub stree<br>
36             <svg fill="none" stroke="currentColor" stroke-linecap="round" stroke-linejoin="round" stroke-width="2" class="w-4 h-4 ml-2">view
37           </p>
38         </div>
39         <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
40
41           <div>
42             <button class="text-white bg-indigo-500 border-0 py-2 px-8 focus:outline-none hover:bg-indigo-600 rounded text-lg">WHAT IS L</button>
43           </div>
44           <p>class="text-base leading-relaxed mt-2">Swag shoidvigolitch literally meditation subway tile tumblr cold-pressed. Gastropub stree<br>
45             <svg fill="none" stroke="currentColor" stroke-linecap="round" stroke-linejoin="round" stroke-width="2" class="w-4 h-4 ml-2">view
46           </p>
47         </div>
48       </div>
49     </div>
50   </div>
51
52   <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
53
54     <div>
55       <div>
56         <div>
57           <div>
58             <div>
59               <div>
60                 <div>
61                   <div>
62                     <div>
63                       <div>
64                         <div>
65                           <div>
66                             <div>
67                               <div>
68                                 <div>
69                                   <div>
70                                     <div>
71                                       <div>
72                                         <div>
73                                           <div>
74                                             <div>
75                                             </div>
76                                           </div>
77                                         </div>
78                                       </div>
79                                     </div>
80                                   </div>
81                                 </div>
82                               </div>
83                             </div>
84                           </div>
85                         </div>
86                       </div>
87                     </div>
88                   </div>
89                 </div>
90               </div>
91             </div>
92           </div>
93         </div>
94       </div>
95     </div>
96   </div>
97
98   <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
99
100    <div>
101      <div>
102        <div>
103          <div>
104            <div>
105              <div>
106                <div>
107                  <div>
108                    <div>
109                      <div>
110                        <div>
111                          <div>
112                            <div>
113                              <div>
114                                <div>
115                                  <div>
116                                    <div>
117                                      <div>
118                                        <div>
119                                          <div>
120                                            <div>
121                                            </div>
122                                          </div>
123                                        </div>
124                                      </div>
125                                    </div>
126                                  </div>
127                                </div>
128                              </div>
129                            </div>
130                          </div>
131                        </div>
132                      </div>
133                    </div>
134                  </div>
135                </div>
136              </div>
137            </div>
138          </div>
139        </div>
140      </div>
141    </div>
142  </div>
143
144  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
145
146    <div>
147      <div>
148        <div>
149          <div>
150            <div>
151              <div>
152                <div>
153                  <div>
154                    <div>
155                      <div>
156                        <div>
157                          <div>
158                            <div>
159                              <div>
160                                <div>
161                                  <div>
162                                    <div>
163                                      <div>
164                                        <div>
165                                          <div>
166                                            <div>
167                                            </div>
168                                          </div>
169                                        </div>
170                                      </div>
171                                    </div>
172                                  </div>
173                                </div>
174                              </div>
175                            </div>
176                          </div>
177                        </div>
178                      </div>
179                    </div>
180                  </div>
181                </div>
182              </div>
183            </div>
184          </div>
185        </div>
186      </div>
187    </div>
188  </div>
189
190  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
191
192    <div>
193      <div>
194        <div>
195          <div>
196            <div>
197              <div>
198                <div>
199                  <div>
200                    <div>
201                      <div>
202                        <div>
203                          <div>
204                            <div>
205                              <div>
206                                <div>
207                                  <div>
208                                    <div>
209                                      <div>
210                                        <div>
211                                          <div>
212                                            <div>
213                                            </div>
214                                          </div>
215                                        </div>
216                                      </div>
217                                    </div>
218                                  </div>
219                                </div>
220                              </div>
221                            </div>
222                          </div>
223                        </div>
224                      </div>
225                    </div>
226                  </div>
227                </div>
228              </div>
229            </div>
230          </div>
231        </div>
232      </div>
233    </div>
234  </div>
235
236  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
237
238    <div>
239      <div>
240        <div>
241          <div>
242            <div>
243              <div>
244                <div>
245                  <div>
246                    <div>
247                      <div>
248                        <div>
249                          <div>
250                            <div>
251                              <div>
252                                <div>
253                                  <div>
254                                    <div>
255                                      <div>
256                                        <div>
257                                          <div>
258                                            <div>
259                                            </div>
260                                          </div>
261                                        </div>
262                                      </div>
263                                    </div>
264                                  </div>
265                                </div>
266                              </div>
267                            </div>
268                          </div>
269                        </div>
270                      </div>
271                    </div>
272                  </div>
273                </div>
274              </div>
275            </div>
276          </div>
277        </div>
278      </div>
279    </div>
280  </div>
281
282  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
283
284    <div>
285      <div>
286        <div>
287          <div>
288            <div>
289              <div>
290                <div>
291                  <div>
292                    <div>
293                      <div>
294                        <div>
295                          <div>
296                            <div>
297                              <div>
298                                <div>
299                                  <div>
300                                    <div>
301                                      <div>
302                                        <div>
303                                          <div>
304                                            <div>
305                                            </div>
306                                          </div>
307                                        </div>
308                                      </div>
309                                    </div>
310                                  </div>
311                                </div>
312                              </div>
313                            </div>
314                          </div>
315                        </div>
316                      </div>
317                    </div>
318                  </div>
319                </div>
320              </div>
321            </div>
322          </div>
323        </div>
324      </div>
325    </div>
326  </div>
327
328  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
329
330    <div>
331      <div>
332        <div>
333          <div>
334            <div>
335              <div>
336                <div>
337                  <div>
338                    <div>
339                      <div>
340                        <div>
341                          <div>
342                            <div>
343                              <div>
344                                <div>
345                                  <div>
346                                    <div>
347                                      <div>
348                                        <div>
349                                          <div>
350                                            <div>
351                                            </div>
352                                          </div>
353                                        </div>
354                                      </div>
355                                    </div>
356                                  </div>
357                                </div>
358                              </div>
359                            </div>
360                          </div>
361                        </div>
362                      </div>
363                    </div>
364                  </div>
365                </div>
366              </div>
367            </div>
368          </div>
369        </div>
370      </div>
371    </div>
372  </div>
373
374  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
375
376    <div>
377      <div>
378        <div>
379          <div>
380            <div>
381              <div>
382                <div>
383                  <div>
384                    <div>
385                      <div>
386                        <div>
387                          <div>
388                            <div>
389                              <div>
390                                <div>
391                                  <div>
392                                    <div>
393                                      <div>
394                                        <div>
395                                          <div>
396                                            <div>
397                                            </div>
398                                          </div>
399                                        </div>
400                                      </div>
401                                    </div>
402                                  </div>
403                                </div>
404                              </div>
405                            </div>
406                          </div>
407                        </div>
408                      </div>
409                    </div>
410                  </div>
411                </div>
412              </div>
413            </div>
414          </div>
415        </div>
416      </div>
417    </div>
418  </div>
419
420  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
421
422    <div>
423      <div>
424        <div>
425          <div>
426            <div>
427              <div>
428                <div>
429                  <div>
430                    <div>
431                      <div>
432                        <div>
433                          <div>
434                            <div>
435                              <div>
436                                <div>
437                                  <div>
438                                    <div>
439                                      <div>
440                                        <div>
441                                          <div>
442                                            <div>
443                                            </div>
444                                          </div>
445                                        </div>
446                                      </div>
447                                    </div>
448                                  </div>
449                                </div>
450                              </div>
451                            </div>
452                          </div>
453                        </div>
454                      </div>
455                    </div>
456                  </div>
457                </div>
458              </div>
459            </div>
460          </div>
461        </div>
462      </div>
463    </div>
464  </div>
465
466  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
467
468    <div>
469      <div>
470        <div>
471          <div>
472            <div>
473              <div>
474                <div>
475                  <div>
476                    <div>
477                      <div>
478                        <div>
479                          <div>
480                            <div>
481                              <div>
482                                <div>
483                                  <div>
484                                    <div>
485                                      <div>
486                                        <div>
487                                          <div>
488                                            <div>
489                                            </div>
490                                          </div>
491                                        </div>
492                                      </div>
493                                    </div>
494                                  </div>
495                                </div>
496                              </div>
497                            </div>
498                          </div>
499                        </div>
500                      </div>
501                    </div>
502                  </div>
503                </div>
504              </div>
505            </div>
506          </div>
507        </div>
508      </div>
509    </div>
510  </div>
511
512  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
513
514    <div>
515      <div>
516        <div>
517          <div>
518            <div>
519              <div>
520                <div>
521                  <div>
522                    <div>
523                      <div>
524                        <div>
525                          <div>
526                            <div>
527                              <div>
528                                <div>
529                                  <div>
530                                    <div>
531                                      <div>
532                                        <div>
533                                          <div>
534                                            <div>
535                                            </div>
536                                          </div>
537                                        </div>
538                                      </div>
539                                    </div>
540                                  </div>
541                                </div>
542                              </div>
543                            </div>
544                          </div>
545                        </div>
546                      </div>
547                    </div>
548                  </div>
549                </div>
550              </div>
551            </div>
552          </div>
553        </div>
554      </div>
555    </div>
556  </div>
557
558  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
559
560    <div>
561      <div>
562        <div>
563          <div>
564            <div>
565              <div>
566                <div>
567                  <div>
568                    <div>
569                      <div>
570                        <div>
571                          <div>
572                            <div>
573                              <div>
574                                <div>
575                                  <div>
576                                    <div>
577                                      <div>
578                                        <div>
579                                          <div>
580                                            <div>
581                                            </div>
582                                          </div>
583                                        </div>
584                                      </div>
585                                    </div>
586                                  </div>
587                                </div>
588                              </div>
589                            </div>
590                          </div>
591                        </div>
592                      </div>
593                    </div>
594                  </div>
595                </div>
596              </div>
597            </div>
598          </div>
599        </div>
600      </div>
601    </div>
602  </div>
603
604  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
605
606    <div>
607      <div>
608        <div>
609          <div>
610            <div>
611              <div>
612                <div>
613                  <div>
614                    <div>
615                      <div>
616                        <div>
617                          <div>
618                            <div>
619                              <div>
620                                <div>
621                                  <div>
622                                    <div>
623                                      <div>
624                                        <div>
625                                          <div>
626                                            <div>
627                                            </div>
628                                          </div>
629                                        </div>
630                                      </div>
631                                    </div>
632                                  </div>
633                                </div>
634                              </div>
635                            </div>
636                          </div>
637                        </div>
638                      </div>
639                    </div>
640                  </div>
641                </div>
642              </div>
643            </div>
644          </div>
645        </div>
646      </div>
647    </div>
648  </div>
649
650  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
651
652    <div>
653      <div>
654        <div>
655          <div>
656            <div>
657              <div>
658                <div>
659                  <div>
660                    <div>
661                      <div>
662                        <div>
663                          <div>
664                            <div>
665                              <div>
666                                <div>
667                                  <div>
668                                    <div>
669                                      <div>
670                                        <div>
671                                          <div>
672                                            <div>
673                                            </div>
674                                          </div>
675                                        </div>
676                                      </div>
677                                    </div>
678                                  </div>
679                                </div>
680                              </div>
681                            </div>
682                          </div>
683                        </div>
684                      </div>
685                    </div>
686                  </div>
687                </div>
688              </div>
689            </div>
690          </div>
691        </div>
692      </div>
693    </div>
694  </div>
695
696  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
697
698    <div>
699      <div>
700        <div>
701          <div>
702            <div>
703              <div>
704                <div>
705                  <div>
706                    <div>
707                      <div>
708                        <div>
709                          <div>
710                            <div>
711                              <div>
712                                <div>
713                                  <div>
714                                    <div>
715                                      <div>
716                                        <div>
717                                          <div>
718                                            <div>
719                                            </div>
720                                          </div>
721                                        </div>
722                                      </div>
723                                    </div>
724                                  </div>
725                                </div>
726                              </div>
727                            </div>
728                          </div>
729                        </div>
730                      </div>
731                    </div>
732                  </div>
733                </div>
734              </div>
735            </div>
736          </div>
737        </div>
738      </div>
739    </div>
740  </div>
741
742  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
743
744    <div>
745      <div>
746        <div>
747          <div>
748            <div>
749              <div>
750                <div>
751                  <div>
752                    <div>
753                      <div>
754                        <div>
755                          <div>
756                            <div>
757                              <div>
758                                <div>
759                                  <div>
760                                    <div>
761                                      <div>
762                                        <div>
763                                          <div>
764                                            <div>
765                                            </div>
766                                          </div>
767                                        </div>
768                                      </div>
769                                    </div>
770                                  </div>
771                                </div>
772                              </div>
773                            </div>
774                          </div>
775                        </div>
776                      </div>
777                    </div>
778                  </div>
779                </div>
780              </div>
781            </div>
782          </div>
783        </div>
784      </div>
785    </div>
786  </div>
787
788  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
789
790    <div>
791      <div>
792        <div>
793          <div>
794            <div>
795              <div>
796                <div>
797                  <div>
798                    <div>
799                      <div>
800                        <div>
801                          <div>
802                            <div>
803                              <div>
804                                <div>
805                                  <div>
806                                    <div>
807                                      <div>
808                                        <div>
809                                          <div>
810                                            <div>
811                                            </div>
812                                          </div>
813                                        </div>
814                                      </div>
815                                    </div>
816                                  </div>
817                                </div>
818                              </div>
819                            </div>
820                          </div>
821                        </div>
822                      </div>
823                    </div>
824                  </div>
825                </div>
826              </div>
827            </div>
828          </div>
829        </div>
830      </div>
831    </div>
832  </div>
833
834  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
835
836    <div>
837      <div>
838        <div>
839          <div>
840            <div>
841              <div>
842                <div>
843                  <div>
844                    <div>
845                      <div>
846                        <div>
847                          <div>
848                            <div>
849                              <div>
850                                <div>
851                                  <div>
852                                    <div>
853                                      <div>
854                                        <div>
855                                          <div>
856                                            <div>
857                                            </div>
858                                          </div>
859                                        </div>
860                                      </div>
861                                    </div>
862                                  </div>
863                                </div>
864                              </div>
865                            </div>
866                          </div>
867                        </div>
868                      </div>
869                    </div>
870                  </div>
871                </div>
872              </div>
873            </div>
874          </div>
875        </div>
876      </div>
877    </div>
878  </div>
879
880  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
881
882    <div>
883      <div>
884        <div>
885          <div>
886            <div>
887              <div>
888                <div>
889                  <div>
890                    <div>
891                      <div>
892                        <div>
893                          <div>
894                            <div>
895                              <div>
896                                <div>
897                                  <div>
898                                    <div>
899                                      <div>
900                                        <div>
901                                          <div>
902                                            <div>
903                                            </div>
904                                          </div>
905                                        </div>
906                                      </div>
907                                    </div>
908                                  </div>
909                                </div>
910                              </div>
911                            </div>
912                          </div>
913                        </div>
914                      </div>
915                    </div>
916                  </div>
917                </div>
918              </div>
919            </div>
920          </div>
921        </div>
922      </div>
923    </div>
924  </div>
925
926  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
927
928    <div>
929      <div>
930        <div>
931          <div>
932            <div>
933              <div>
934                <div>
935                  <div>
936                    <div>
937                      <div>
938                        <div>
939                          <div>
940                            <div>
941                              <div>
942                                <div>
943                                  <div>
944                                    <div>
945                                      <div>
946                                        <div>
947                                          <div>
948                                            <div>
949                                            </div>
950                                          </div>
951                                        </div>
952                                      </div>
953                                    </div>
954                                  </div>
955                                </div>
956                              </div>
957                            </div>
958                          </div>
959                        </div>
960                      </div>
961                    </div>
962                  </div>
963                </div>
964              </div>
965            </div>
966          </div>
967        </div>
968      </div>
969    </div>
970  </div>
971
972  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
973
974    <div>
975      <div>
976        <div>
977          <div>
978            <div>
979              <div>
980                <div>
981                  <div>
982                    <div>
983                      <div>
984                        <div>
985                          <div>
986                            <div>
987                              <div>
988                                <div>
989                                  <div>
990                                    <div>
991                                      <div>
992                                        <div>
993                                          <div>
994                                            <div>
995                                            </div>
996                                          </div>
997                                        </div>
998                                      </div>
999                                    </div>
1000                                  </div>
1001                                </div>
1002                              </div>
1003                            </div>
1004                          </div>
1005                        </div>
1006                      </div>
1007                    </div>
1008                  </div>
1009                </div>
1010              </div>
1011            </div>
1012          </div>
1013        </div>
1014      </div>
1015    </div>
1016  </div>
1017
1018  <div class="p-4 md:w-1/3 sm:mb-0 mb-6">
1019
1020    <div>
1021      <div>
1022        <div>
1023          <div>
1024            <div>
1025              <div>
1026                <div>
1027                  <div>
1028                    <div>
1029                      <div>
1030                        <div>
1031                          <div>
1032                            <div>
1033                              <div>
1034                                <div>
1035                                  <div>
1036                                    <div>
1037                                      <div>
1038                                        <div>
1039                                          <div>
1040                                            <div>
1041                                            </div>
1042                                          </div>
1043                                        </div>
1044                                      </div>
1045                                    </div>
1046                                  </div>
1047                                </div>
1048                              </div>
1049                            </div>
1050                          </div>

```


BIBLIOGRAPHY

- <https://smartbridge.teachable.com/courses/1450164/lectures>
- <https://www.math.ucla.edu/~minchen/doc/ImgCapGen.pdf>
- <https://realpython.com/python-web-applications/>