

000
001
002
003
004
005
006
007
008
009
010
011

Abstract

With this project, I have developed a methodology using which people can identify the names of the books and also, find out dimensional information regarding the text book. The idea is implemented such that users will turn on the live feed using their cameras and the program while working in the background will recognize the key features in the frame to classify the book correctly. It is important to gather the book features before a certain book is tested, because similar to image recognition, we would not be able to determine the name of the book if it is not in the database. Techniques such as image warping, knn (k-nearest neighbors) matching for classification, and contour parameters are used to precisely measure the dimensions of the book. Some of the challenges such as depth perception are also tackled in this project by utilizing the geometric manipulation and by creating an environment that best identifies the framework. An advantage of this technique is that we can store the features extracted from the database and in the future, we may create eigen faces to reduce the required memory size and it will allow for faster classification of objects in an image/frame. I also compared the book recognition task by implementing Convolutional Neural Network (CNN) and Artificial Neural Network (ANN) networks to analyze the accuracy and loss metrics.

1. Introduction

With online shopping reaching new heights, it is inevitable that most companies would require their products and services to be offered online. As can be witnessed, the usual goal of these companies is not just to bring their products online but also to add an additional feature which allows their users to compare different products on their website and also against their competitor's products. In order to achieve these objectives, it is logical to invest in the development of techniques and algorithms that can quickly and accurately identify the products. These trends are also seen in the realm of book selling market and therefore, it is important that developers create book cover recognition

Book Cover Recognition

COMPENG 4TN4: Image Processing - Project 2 Report

Harneet Singh - 400110275

algorithms. The aim of these algorithm should be to classify the books based on genre, author, price, name and other relevant information related to selling points of books. The goal of my project is to classify the images based on the features associated with the book cover. In the future, more work may be performed to create apps and websites that can offer these services online and in real time to better assist the users.

In this project, my program extracts features from images containing books and creates descriptors for the said features. The goal of adding descriptors is to classify the image and draw a bounding box around the book in a real time camera feed (or video). Before finding features in the image, first job would be to create a database that will contain images with their labels. In order to create a strong database, image pre-processing steps will be applied to remove unwanted variations from the data images. After resizing the images in the database and also from the camera feed, we can compare the descriptors between two images using knn matching. I have tested the classification through knn and FLANN matchers. FLANN comprises algorithms that are optimized for fast nearest neighbor search in large datasets and for high dimensional features [5]. In the end, object needs to be localized in each frame of the video and a marker will be attached to it to specify its class. This way object recognition task will be achieved for this project.

Additionally, a stretch goal is implemented to measure the dimensions of the book and add a feature that calculates total pages (approx.) when the book thickness is in the frame. To achieve this functionality, I had to apply filters such as resizing, smoothing filter, edge detection to binarize the image and morphological operators to get the best input image data.

Also, after implementing the object recognition steps, I have attempted to compare the results of my project with pre-existing object recognition classifiers.

1.1. Related Work

At Stanford University, Linfeng Yang and Xinyu Shen created a GUI in MATLAB to determine information about a desired book in real time [7]. After taking pictures of

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

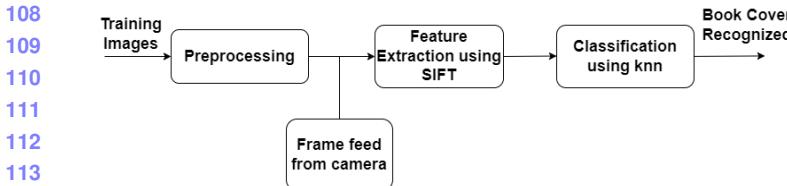


Figure 1. Proposed Method Flowchart

the book cover, the app detected feature of the input image based on Maximally Stable Extremal Region (MSER) algorithm. They captured the textual and non-textual information separately based on morphological difference. To get optimal results, they created a text character alignment algorithm which improves the accuracy of original text detection. Later, they compared their algorithm with MATLAB's built-in OCR algorithm and a commonly used open source OCR. Further, to eliminate false detection inhibition and word correction, they implemented a post detection algorithm and natural language processing. On the online platform, their algorithm attained accuracy numbers of greater than 86%.

In the past, work has been conducted on OCR technology namely, Tesseract OCR engine [6]. OCR engine has been used to perform identification task on natural language. In their experiments, they analyzed the Tesseract OCR along with Transym and their conclusion was that Tesseract OCR is more stable when it comes to large image databases.

Previously, other techniques have been used to extract features from images such as SIFT algorithm. Technically speaking, in my opinion, any work performed to improve the accuracy of object classification in an image task is relevant to the book cover recognition because similar principles can be applied in book cover recognition. The simplicity and efficacy of SIFT algorithm inspired me to work on this project using SIFT algorithm while performing preprocessing steps. The biggest advantage of applying SIFT algorithm is that it can identify the features without regard to orientation and scale (to a certain degree) of the book.

2. Proposed Method with Results

First step is create a database with train images and their labels. After that, preprocessing steps have been applied to acquire consistent images and images are converted to grayscale to reduce the number of channels associated with train images. To extract the features, descriptors and keypoints are stored for each image using SIFT algorithm and subsequently, knn-matcher is used to classify the camera frames accurately. Finally for measuring the book dimensions, preprocessing steps are applied which will be explained in the following sections. Flowchart is shown in Fig. 1.



Figure 2. Collage of Train Images (Original)



Figure 3. Collage of Train Images (Resized and grayscaled)

2.1. Train Images Dataset

To create a dataset, I am utilizing the Harry Potter book collection. Book cover images are found using an online source and labels are attached to each image as the image's name (as shown in Fig. 2).

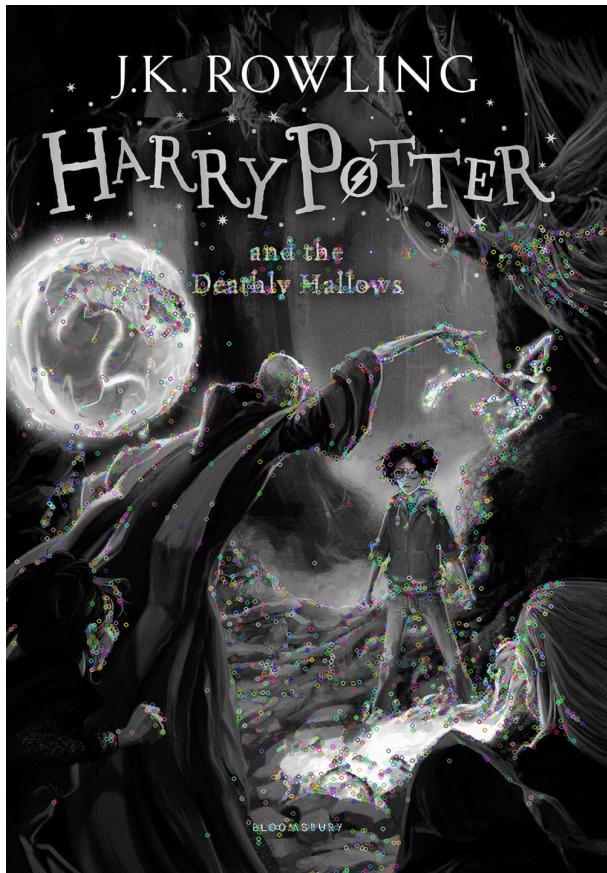
Labels: ['HP and the Chamber of Secrets', 'HP and the Deathly Hallows', 'HP and the Goblet of Fire', 'HP and the Half-Blood Prince', 'HP and the Order of the Phoenix', "HP and the Philosopher's Stone", 'HP and the Prisoner of Azkaban']

2.2. Resizing and Grayscale

After gathering the images, their sizes are reduced to a scale of 60% (empirically tested, program allows for alteration of scale factor) and images are converted to grayscale (as shown in Fig. 3).

108	162
109	163
110	164
111	165
112	166
113	167
114	168
115	169
116	170
117	171
118	172
119	173
120	174
121	175
122	176
123	177
124	178
125	179
126	180
127	181
128	182
129	183
130	184
131	185
132	186
133	187
134	188
135	189
136	190
137	191
138	192
139	193
140	194
141	195
142	196
143	197
144	198
145	199
146	200
147	201
148	202
149	203
150	204
151	205
152	206
153	207
154	208
155	209
156	210
157	211
158	212
159	213
160	214
161	215

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243



244 Figure 4. Extracted descriptors and keypoints of a train image
245 using SIFT algorithm

246 2.3. Extracting features from train images

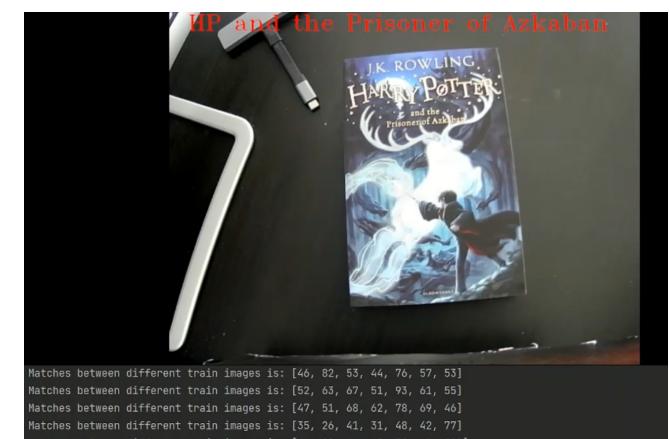
247 Using SIFT algorithm to find the keypoints and descriptors
248 of each train image and storing it in separate lists.
249 These keypoints and descriptors will be used later on for
250 comparison with the camera frames. An example of identified
251 descriptors is shown in Fig. 4 using a train image. Similar
252 results are obtained for other train images as well (please
253 view the attached results).

254 2.4. Comparison of train images features against 255 camera-frame descriptors

256 Camera is turned on and each frame is used to detect the
257 features (using SIFT) within the frame. These descriptors
258 are compared with the train images feature dataset. Brute
259 force matcher is used to get k (=2) best matches between
260 the frame and train images. Knn technique is deployed to
261 find matches (cross-validated) and ratio test is applied. An
262 example of matches between camera-frame and train image
263 is shown in Fig. 5.



264 Figure 5. Matches between camera-frame and a train image



265 Matches between different train images is: [46, 82, 53, 44, 76, 57, 53]
266 Matches between different train images is: [52, 65, 67, 51, 93, 61, 55]
267 Matches between different train images is: [47, 51, 68, 62, 78, 69, 46]
268 Matches between different train images is: [35, 26, 41, 31, 48, 42, 77]
269 Matches between different train images is: [44, 61, 40, 69, 78, 40, 253]
270 Matches between different train images is: [40, 44, 39, 38, 59, 35, 257]

271 Figure 6. Result: Identified Book Cover

272 2.5. Thresholding to attain the best match

273 Thresholding technique is used to eliminate the train images
274 with low matches. After finding the good matches between
275 the train images and camera-feed, a counter is used to
276 determine the number of good matches of each train image.
277 The image with count value above the threshold and with
278 the highest count is determined to be a good match for the
279 book cover. An example of this result is shown in Fig. 6.
280 (Please watch the video in Step 5 to find other results)

281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323

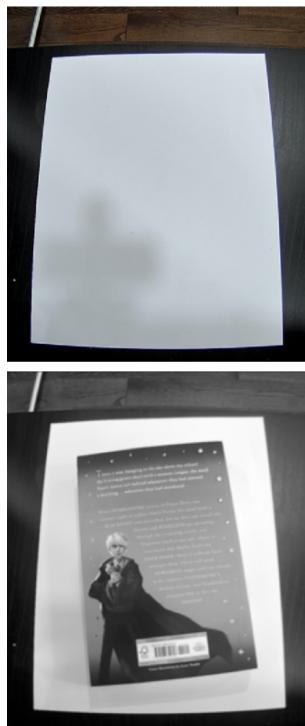
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350

Figure 7. Gaussian Filtered Frame

351
352
353

Determining book dimensions and number of pages method

354
355
356

Following sections explain how the book dimensions and number of pages were determined.

357
358

2.3-b Smoothing Filter

359
360
361

After the steps of 2.2. Resizing and Grayscale, gaussian filter of kernel size equal to 5 is applied to remove noise. Result is shown for a frame in Fig. 7.

362

2.4-b Binarizing Frame

363
364

Canny edge detector is applied to binarize the frame image. To determine optimal Canny edge detectors threshold values, sigma value and median of the image are used and thereby, low and high threshold are calculated. The idea is to pick values that are close to the median. Result is shown for a frame in Fig. 8.

371

2.5-b Morphological Operations

372
373
374
375
376
377

Morphological operators are applied to remove unwanted blobs and ultimately, clean up the threshold image from previous step. Iterations of dilate and erode are performed and result is shown for a frame in Fig. 9.



Figure 8. Edges Detected in the frame



Figure 9. Morphological Operations on threshold image

2.6-b Finding external contours and bounding box

To be able to figure out the size of the book, we need an object in the frame that will allow us to determine the relationship between number of pixels and distance in 'millimeter' in the frame. To achieve this, I am using an A4 paper (size 210mm x 297mm), and book is placed on this paper to correctly measure the size. Area parameter is used to find the contour in the frame. For this setup to work, we need to ensure that the depth is correctly perceived by the camera and my approach to tackle this issue will be explained later. Contour and bounding box results are shown for a frame in Fig. 10.

2.7-b Warp image

From Fig. 10, we can clearly see that that camera is not aligned properly. Based on the bounding box, we can state that through the camera's eye the width and height is not consistent. To fix this issue, we need to warp the image by calculating transformation matrix and applying warp perspective transformation on the image. Result is shown for a frame in Fig. 11.

2.8-b Internal Contour

Now, we can find the internal contour i.e., the book on the A4 page. Result is shown for a frame in Fig. 12.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

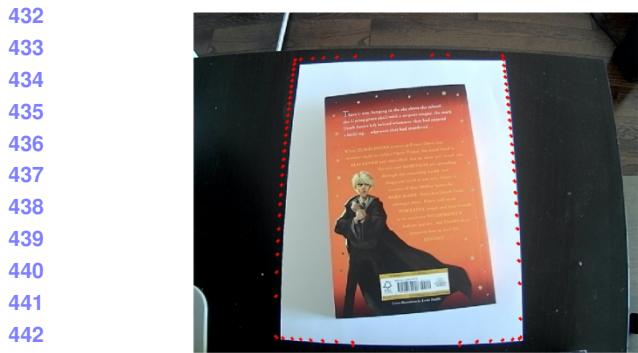


Figure 10. Contour and Bounding Box of A4 paper



Figure 11. Warped Perspective of A4 paper

2.9-b Determine all three dimensions of the book and compute number of pages

To find the height and width of the image, we can use the height and width parameters of the bounding box of the contour. Result is shown in Fig. 13.

However, to measure the thickness of the book, we need to account for the elevation of the book contour. Since the dimensions of the object changes as we move the object closer to the camera, we need to account for the elevation bias when the book is not laying flat as shown in Fig. 12. With the current setting, we can precisely measure the width

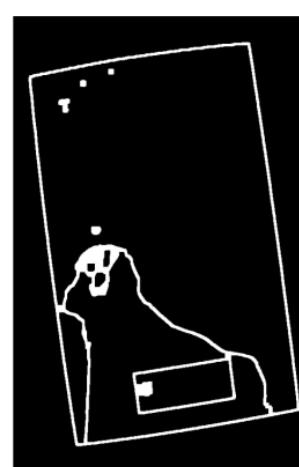


Figure 12. Book Contour on A4 paper

and height of the book as shown in Fig. 12, because the elevation in the contour is negligible. In this case, the thickness of the book does not play a significant role in skewing the depth perception of the camera. But when we place the book along its length and thickness against the A4 paper (as shown in Fig. 14), the width of the book produces an elevation bias. Measuring dimensions of an object requires geometric manipulation as explained below.

As explained earlier, in Fig. 14 the contour is lifted higher and the contour seen by the camera has a larger area than the actual area of the book on the A4 paper. To solve this problem, we need to analyze Fig. 15 and we can deduce the following: (let's say $height/2 = h/2$)

$$\frac{x}{d} = \tan\beta$$

$$\frac{h/2}{d + w} = \tan\alpha$$

For my case, $d = 31"$ and $(d + w) = 36"$. In reality, $h/2$ is equal to x , therefore,

$$31.\tan\beta = 36.\tan\alpha$$

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

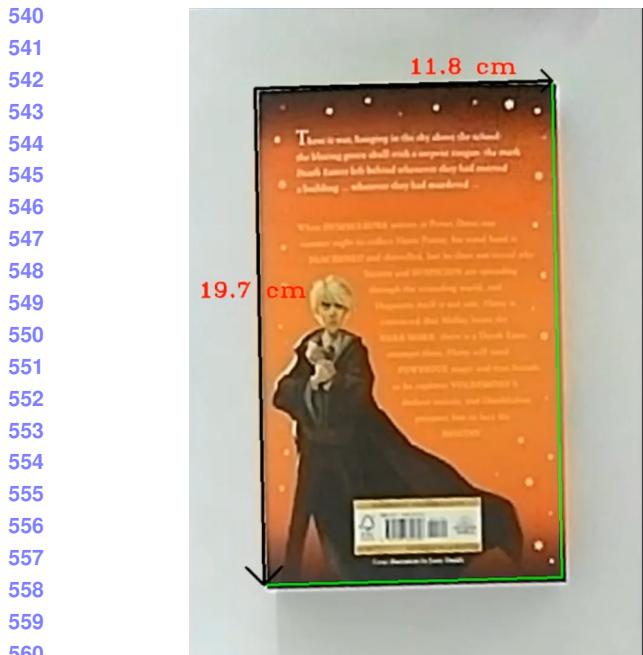


Figure 13. Result: Height and Width of the book



Figure 14. Height and Thickness side of the book

$$\frac{\tan \alpha}{\tan \beta} = 0.86 = 0.9(\text{approx.})$$

Hence, I am using an elevation bias of 0.9 to correct the final output, because as 'd' grows, $\frac{d}{d+w}$ approaches value of 1. We can imagine 'x' is calculated as a large number and therefore, we need to reduce it by multiplying by elevation bias of 0.9. Final, result is shown along with the number of pages of the book in Fig. 16. To find number of pages, I determined the ratio of number of pixels along the height and used this ratio to find the number of pages per millimeter. To visualize the results to further details, please watch the videos attached with the report on Avenue to learn.

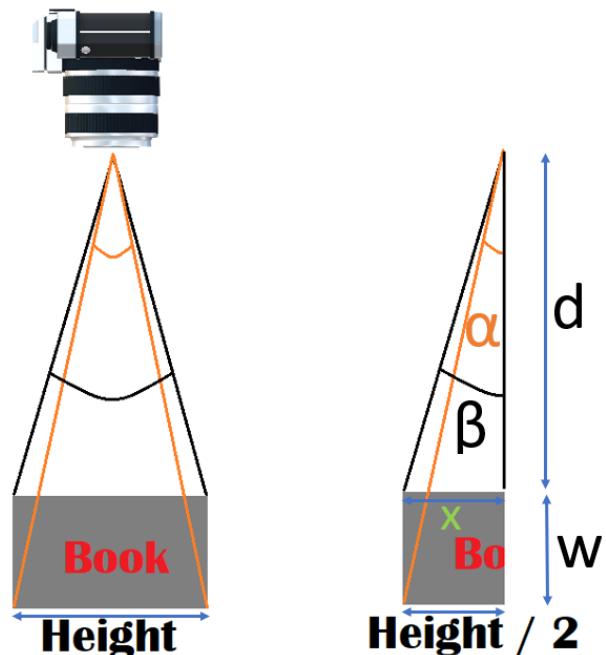


Figure 15. Camera Depth Perception

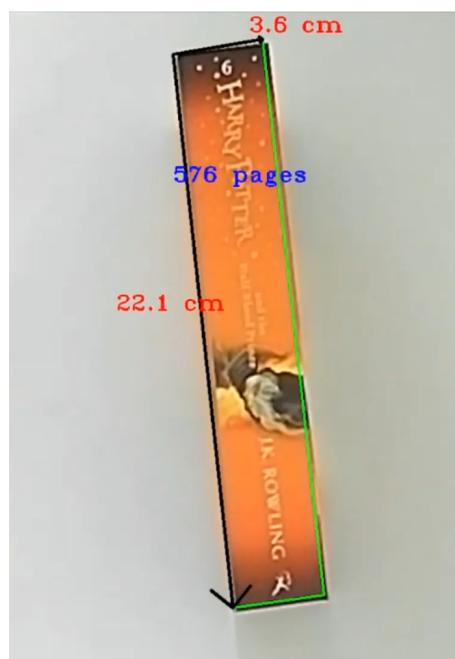


Figure 16. Result: Measured thickness, height, and number of pages

3. Experimental Results for Comparison with Proposed Method

To compare my method, I have implemented book cover recognition task using ANN and CNN. Results of these will

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

```

648 Model: "sequential"
649
650 Layer (type)          Output Shape         Param #
651 =====
652 flatten (Flatten)     (None, 2500)        0
653 dense (Dense)         (None, 3000)        7503000
654 dense_1 (Dense)       (None, 1000)        3001000
655 dense_2 (Dense)       (None, 500)         500500
656
657 Total params: 11,004,500
658 Trainable params: 11,004,500
659 Non-trainable params: 0

```

Figure 17. ANN Model Summary

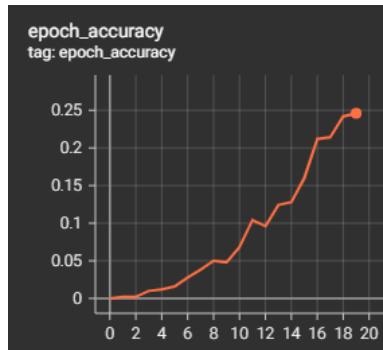


Figure 18. ANN Model Accuracy with 20 epochs

be presented in the following sections. Note that these experiments were conducted to understand the concepts behind the implementation, process of implementation and to compare the accuracy between different methods.

Book cover dataset has been collected from Kaggle [1] and only a portion of the dataset was used to compare the performance (500 train images) to quickly train the models.

3.1. Artificial Neural Network (ANN)

Fig. 17 shows the model that was created to test the accuracy and loss metrics of ANN. Only four layers are used with labels represented as sparse categorical.

Results of accuracy and loss can be visualized in Fig. 18 and Fig. 24 respectively. After 20 epochs, the model reached accuracy of 25.6% with loss of 5.4649 and while testing the model with test images, the model performed at accuracy rate of 45.0%. Similar images were used to test the trained model, however another challenge for the model could be developed by altering the scale and orientation of the test images. Result of test image is shown in Fig. 20 and Fig. 21, in this case model classified correctly and incorrectly respectively.

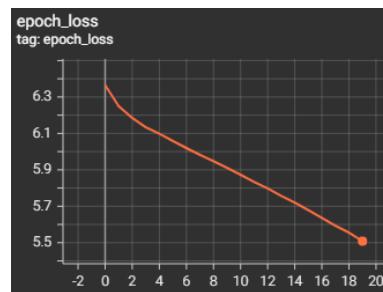


Figure 19. ANN Model Loss with 20 epochs

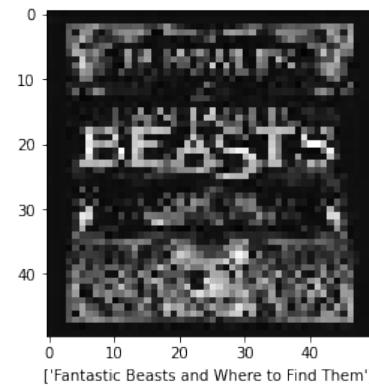


Figure 20. Result: Correct - ANN Test Classification

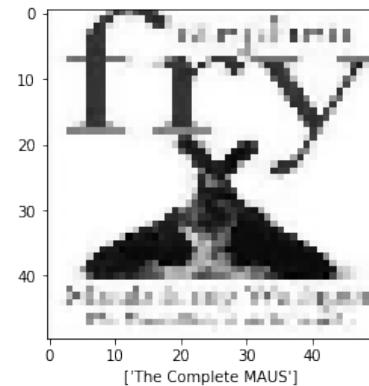


Figure 21. Result: Incorrect - ANN Test Classification

3.2. Convolutional Neural Network (CNN)

Fig. 22 shows the model that was created to test the accuracy and loss metrics of CNN. Only seven layers consisting of convolution, pooling, flatten and dense layers are used with labels represented as sparse categorical.

Results of accuracy and loss can be visualized in Fig. 23 and Fig. 24 respectively. After 20 epochs, the model reached accuracy of 88.8% with loss of 0.4668 and while testing the model with test images, the model performed at

```

756
757 Model: "sequential_1"
758
759 Layer (type) Output Shape Param #
760 =====
761 conv2d (Conv2D) (None, 48, 48, 32) 320
762 max_pooling2d (MaxPooling2D (None, 24, 24, 32) 0
763 )
764 conv2d_1 (Conv2D) (None, 22, 22, 64) 18496
765 max_pooling2d_1 (MaxPooling2D (None, 11, 11, 64) 0
766
767 flatten_1 (Flatten) (None, 7744) 0
768 dense_3 (Dense) (None, 64) 495680
769 dense_4 (Dense) (None, 500) 32500
770
771 Total params: 546,996
772 Trainable params: 546,996
773 Non-trainable params: 0

```

Figure 22. CNN Model Summary

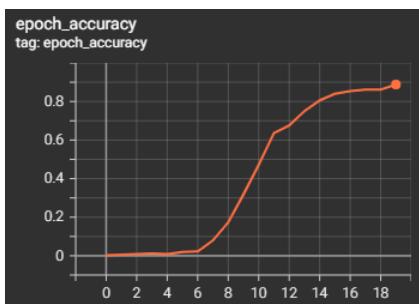


Figure 23. CNN Model Accuracy with 20 epochs

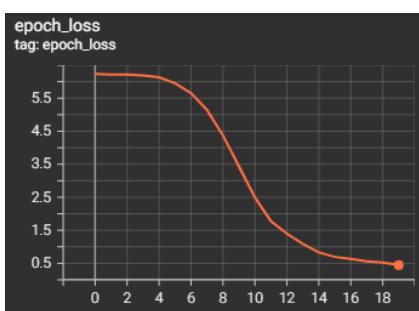


Figure 24. CNN Model Loss with 20 epochs

accuracy rate of 100.0%. Similar images were used to test the trained model, however another challenge for the model could be developed by altering the scale and orientation of the test images. Result of test image is shown in Fig. 25, in this case model classified correctly.

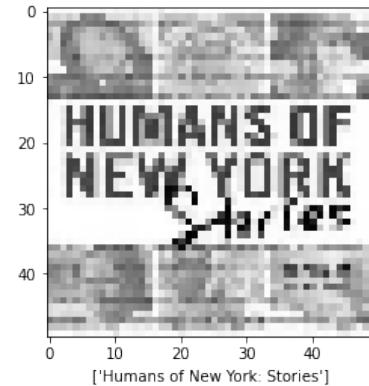


Figure 25. Result: Correct CNN Test Classification

4. Conclusion and Future Work

Utilization of SIFT has proven to be a strong candidate for feature detection in the cover pages of the books. However, similar to training CNN and ANN models, SIFT algorithm requires time to extract the features. An alternative is to extract and store the descriptors and keypoints for the entire dataset. In the future, a fruitful addition would be to incorporate alphanumeric character recognition techniques such as OCR as well. On the other side, as we have seen that CNN can clearly outperform ANN in almost all aspects. However, aside from the time complexity, a drawback of ANN and CNN is that it is not scale and orientation invariant. To improve the identification of book covers regardless of scale and orientation, we need to perform tasks of data augmentation and appropriately train the ANN and CNN models. In this project, my focus was to correctly recognize the book covers and determine the dimensions and number of pages of the book. These requirements were satisfied. In the future, another approach could be to not only link the book cover with its label but also, associate it with price, genre, author and other relevant information.

References

- [1] LUKA ANICIN. Book covers dataset. <https://www.kaggle.com/datasets/lukaaninicin/book-covers-dataset>. 7
- [2] Tsai S. S. Schroth G. Chen D. M. Grzeszczuk R. Girod B. Chen, H. Robust text detection in natural images with edge-enhanced maximally stable extremal regions. (2011, September) In 2011 18th IEEE International Conference on Image Processing (pp. 2609-2612). IEEE.
- [3] codebasics. deep-learning-keras-tf-tutorial. https://github.com/codebasics/deep-learning-keras-tf-tutorial/blob/master/16_cnn_cifar10_small_image_classification/cnn_cifar10_dataset.ipynb.

864	[4] Harrison@pythonprogramming.net. Classifying cats vs dogs	918
865	with a convolutional neural network on kaggle. https://pythonprogramming.net/convolutional-neural-network-kats-vs-dogs-machine-learning-tutorial/ .	919
866		920
867		921
868		922
869	[5] OpenCV-Python. Feature matching. https://opencv24-python-tutorials.readthedocs.io/en/latest/py_tutorials/py_feature2d/py_matcher/py_matcher.html . 1	923
870		924
871		925
872		926
873	[6] R. Smith. An overview of the tesseract ocr engine. (2007). 2	927
874		928
875	[7] Linfeng Yang and Xinyu Shen. Book cover recognition. https://web.stanford.edu/class/ee368/Project_Autumn_1617/Reports/report_yang_shen.pdf . 1	929
876		930
877		931
878		932
879		933
880		934
881		935
882		936
883		937
884		938
885		939
886		940
887		941
888		942
889		943
890		944
891		945
892		946
893		947
894		948
895		949
896		950
897		951
898		952
899		953
900		954
901		955
902		956
903		957
904		958
905		959
906		960
907		961
908		962
909		963
910		964
911		965
912		966
913		967
914		968
915		969
916		970
917		971