

Ground detection: Road Detection

A. Introduction

Road detection is an important task, together with obstacle detection perhaps one of the main tasks in autonomous navigation. In the past, road and obstacle recognition have been heavily dependent on information gained from LIDAR and other active surface recognition methods. However use of LIDAR may pose a health hazard in certain situations. Hence a passive detection method based solely on visual feedback is essential.

We have explored and implemented two methods for road detection. The first method uses optical flow technique, and second one is based on semantic segmentation using Deep Learning. Semantic segmentation is the task to categorize each pixel into a class. Many deep learning networks developed for semantic segmentation can be fine-tuned for road segmentation. We have studied different models and trained two of them to develop a road detection system. We used existing datasets like Cityscapes for training our models. We have used the images clicked around the UMD campus for further testing of the models. The output of this project is a comparison between the performance of selected Deep Learning models/other methods.

Although a driveable path for an autonomous vehicle may not be limited to roads, (for e.g. an autonomous scooter might be able to drive on sidewalks), our project looks at detection of the roads as the main goal.

B. Methods explored

We implemented two methods for road detection. Our initial plan was to implement semantic segmentation using deep learning. But before going ahead with that , we also explored optical flow as a method for detecting.

- Optical Flow
 - Sparse Optical Flow
 - Dense Optical Flow
- Semantic segmentation using deep learning We trained and compared the two models:
 - Unet with Resnet18 backbone
 - DeepLabV3 with Resnet18 backbone

C. Optical Flow

The idea behind using optical flow was that for a given set of images a vehicle moving on a road, the road is assumed to have the least number of features. Hence it will have the lowest velocity (close to zero). Thus we will by using velocity as our metric, we will be able to classify the

environment into different objects and will be able to identify the road out of them.

We tried out two methods int Optical FLow.

C.1. Sparse Optical Flow

Sparse optical flow refers to a technique that estimates the motion of selected key points or features between consecutive frames of an image or video sequence. Sparse optical flow focuses on a subset of interest points, and not all pixel points.

We take the first frame, detect some Shi-Tomasi corner points in it, then we iteratively track those points using Lucas-Kanade optical flow. The results obtained are not satisfactory, since for a case of a car moving on the road, a sparse method yields features and can be more useful to estimate the path of movement or velocity, but inefficient for road detection.

Also, it did not work for the case where the car is moving too fast, the trajectory came out to be not very useful.

C.2. Dense Optical Flow

Unlike sparse optical flow, which focuses on tracking selected key points, dense optical flow calculates motion vectors for every pixel in the image. The objective of dense optical flow is to capture the detailed motion information within an image or video. It provides a dense field of motion vectors that represent the displacement of each pixel between frames. These motion vectors describe the magnitude and direction of pixel movement, allowing for a comprehensive understanding of the dynamics in the scene

One commonly used dense optical flow algorithm is the Gunnar Farneback method, which is used in our case.

For the same set of videos in the previous section :

- Dense optical flow is found using the above algorithm.
- We get a 2-channel array with optical flow vectors (u, v).
- We find their magnitude and direction. We color code the result for better visualization. →Direction corresponds to the Hue value of the image. Magnitude corresponds to the value plane.

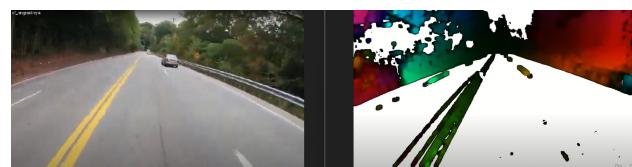


Figure 1. Dense optical flow output (Gunnar - Farneback)

- 108 • Additionally, we changed the dark colors in the color-coded frames to white, for better understanding.
 109
 110
 111 • The algorithm computes the relative velocity of pixels. In the case of road planes, since the road color pretty much remains the same, their velocity vectors
 112 are mostly constant, which is how we detect roads.
 113
 114 • There are a few limitations that come with using such
 115 an approach for road detection. This method is bound
 116 to fail in the case of the absence of a significant change
 117 of pixels from one frame to the other, typically in the
 118 case of a road with barren land (deserted) surroundings.
 119 The magnitude thresholding, in this case, will be
 120 challenging as the algorithm might not be able to de-
 121 tect the pixels corresponding to the surroundings to be
 122 moving, which makes it difficult to tell the road apart.
 123
 124 • Another challenge is high computational complexity.
 125 Due to the high computational complexities of this al-
 126 gorithm, it might not be very feasible to deploy this in
 127 real-time situations, for road detection-related applica-
 128 tions.
 129
 130

D. Semantic segmentation using Deep Learning

131 As all the team members were new to the field of deep
 132 learning, the first step was to understand what semantic seg-
 133 mentation and deep learning is. Some part of it was taught
 134 during the class. We also followed a few tutorials on py-
 135 torch website to learn the practical implementation.
 136

137 To get started, we tried to use some pre-trained models
 138 available on PyTorch hub and tested with some of the out-
 139 door images captured at the UMD campus.
 140

141 We explored and tested the following models:
 142

- 143
 144 • DeepLabV3 with ResNet101 CNN architecture
 145
 146 • FCN(Fully Convolutional Networks) with ResNet101
 147 backbone
 148
 149 • Unet with Resnet18 and Resnet34 backbone
 150

151 We didn't get expected results from the above pre-
 152 trained models.
 153

- 154 • HybridNets HybridNets is a state-of-art perception
 155 network, specially for the drivable area segmentation
 156 and lane detection. It was working well for the images
 157 we tested.
 158
 159 • A benchmark semantic segmentation network from
 160 the ADE20K MIT Scene Parsing Benchmark: This
 161 model uses resnet50dilated and ppm deeplab archi-
 162 tectures for encoder and decoder respectively. It was
 163 trained by authors on ADE20K Dataset which has 150



164
 165
 166
 167
 168
 169
 170
 171
 172
 173
 174
 175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214
 215

Figure 2. HybridNet Results

semantic categories including road, grass etc. We used
 it as a pretrained model. The results of segmentation
 are shown below.



Figure 3. Results1

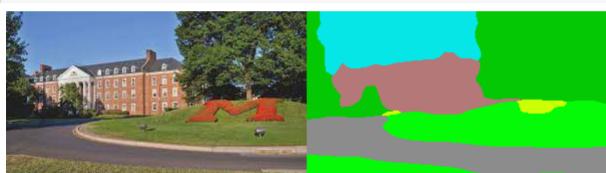


Figure 4. Results2



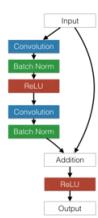
Figure 5. Results3



Figure 6. Results4

216 E. Information about architectures used 270 217

218 E.1. Resnet 271 219



220 Figure 7. ResNet 284

221 ResNet is a short form for Residual Network. The key 285 feature about this architecture is the introduction of 286 residual blocks. It solves the problem of vanishing 287 gradients. The ResNet implements skip connections by 288 creating shortcut that bypass one or more layers in a 289 network, allowing the 290 gradient to flow directly across multiple layers. This 291 property allows resnet to train extremely deep neural 292 networks. 293

294 There are several variations of the ResNet architecture, 295 such as ResNet-18, ResNet-34, ResNet-50, ResNet-101, 296 and ResNet-152, which differ in terms of depth. The 297 number after "ResNet" indicates the number of layers in the 298 network. The deeper variants (e.g., ResNet-101, ResNet-152) 299 have been shown to achieve better performance on certain 300 tasks but require more computational resources. 301

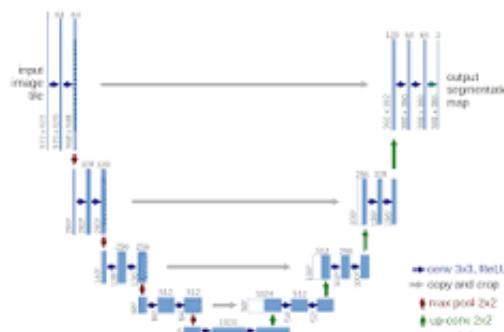
302 ResNet has been widely used and has achieved state-of- 303 the-art results in various computer vision tasks, including 304 image classification, object detection, and semantic 305 segmentation. It has also been a foundation for many subsequent 306 architectures and has greatly influenced the field of 307 deep learning. 308

309 E.2. Unet architecture 330

311 The U-Net architecture is a popular deep learning 312 architecture for semantic segmentation tasks[3].It gets its name 313 from its U-shape structure, which is shown below. It is 314 particularly effective for tasks where precise localization is 315 required, such as segmenting objects in medical images. 316

317 U-Net has encoder-decoder type architecture.The 318 encoder takes images and transforms it into a compact and 319 abstract representation.It consists of a series of convolutional 320 and pooling layers, where each convolutional layer is 321 typically followed by a rectified linear unit (ReLU) activation 322 function. The pooling layers reduce the spatial dimensions 323 of the feature maps while increasing the number of channels. 324

325 The decoder performs upsampling and reconstructs the 326 segmented output based on the features learned in the 327 encode. It implements transpose convolutions for this. This 328 skip connection allows the decoder to access fine-grained 329

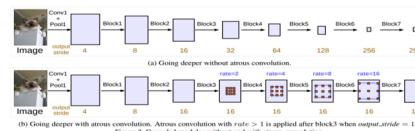


270 Figure 8. UNet Architecture 281

271 information from earlier layers, aiding in precise 272 localization. The final layer of the U-Net architecture is a 1x1 273 convolutional layer that maps the high-dimensional feature 274 maps to the desired number of output channels, representing 275 the segmented output. 276

277 Overall, the U-Net architecture combines the benefits of 278 both convolutional neural networks for feature extraction 279 and fully convolutional networks for dense prediction. The 280 skip connections enable the network to leverage multi-scale 281 features, which is especially useful for accurate segmentation 282 of objects with varying sizes. 283

284 E.3. DeepLabV3 295



296 Figure 9. DeepLabV3 Architecture 304

297 DeepLabv3 is a popular semantic segmentation 305 architecture. It is a fully Convolutional Neural Network (CNN) 306 model to tackle the problem of semantic segmentation. To 307 handle the problem of segmenting objects at multiple scales, 308 modules are designed which employ dilated convolution. It 309 is an end-to-end trainable deep learning system. 310

311 Deeplabv3 cemented its status as one of the go-to 312 models for segmentation because of its high speed,better 313 accuracy, architectural simplicity and generalizability to custom 314 tasks. 315

316 From an architectural standpoint, DeepLabv3 uses the 317 ResNet models (trained on Imagenet) as the backbone along 318 with the widely popular Atrous(dilated) Convolution and 319 Atrous Spatial Pyramid Pooling (ASPP) module. 320

321 F. Dataset 360

322 We used the Cityscape dataset for training our model. 323 The Cityscapes dataset is a widely used benchmark dataset 324

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
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
for urban scene understanding and computer vision tasks, particularly in the context of autonomous driving. It provides a large-scale collection of high-quality images with pixel-level annotations for various semantic segmentation tasks.



Figure 10. Cityscape Dataset

Below are the key features of the Cityscapes dataset:

- 1. Data Collection:** The dataset consists of images captured in different cities across Germany. The images were taken from a car-mounted stereo camera setup, resulting in high-resolution and high-quality RGB images.
- 2. Semantic Segmentation:** It provides pixel-level annotations for 30 different classes, including road, sidewalk, buildings, vehicles, pedestrians, traffic signs, and vegetation. The annotations enable the training and evaluation of models for semantic segmentation tasks in urban environments.
- 3. Fine-annotated Subset:** The dataset also includes a "fine" subset, which has more detailed annotations for selected object classes. This subset contains more accurate and dense annotations for object boundaries and instance-level information.
- 4. Coarse-annotated Subset:** In addition to the fine-annotated subset, the dataset includes a "coarse" subset, which provides annotations with lower spatial resolution. The coarse annotations are intended for faster processing and can be useful for certain applications where fine-grained details are not necessary.
- 5. Additional Data:** Cityscapes also provides complementary data, including disparity maps, which are useful for tasks such as depth estimation, and camera cali-

bration files, which facilitate accurate geometric alignment of the data.

G. Approach

Our approach to implement and train deep learning model for the task of semantic segmentation can be described in the following steps.

1. **Data Preparation:** As mentioned in previous section, We are using Cityscapes dataset from torchvision.datasets for training our models. We are using Fine-annotated subset in particular. It has images and annotations split in train, val and test sub-directories to use it for training, validation and testing purpose. It has 30 classes such as road, sidewalk, vegetation, buildings etc. However, for our application, we just need to detect the road. We want to know whether the pixel is road or not. Therefore, we removed unwanted classes and rectified the labels of wanted classes. We also performed data augmentation using albumentation library. In data augmentation technique, various transformations can be applied on original data to increase the size of training data. It improves the performance of model, and make it robust.

2. **Build and Train Models:** We used PyTorch Lightning for this task. It is a lightweight version of PyTorch that simplifies the training and development of neural networks. It provides a high-level interface and abstractions for common deep learning tasks, allowing users to focus more on the model architecture and experimentation rather than the low level code.

Our model class has been inherited from LightningModule in PyTorch Lightning. We selected Unet and DeeplabV3 architectures. We selected resnet18 encoder initialized with 'imagenet' pre-trained weights for both the models. We have used AdamW Optimiser, which is a variant of the Adam optimizer that introduces weight decay, along with combining adaptive learning rates to accelerate convergence during training. The training loop is handled by Trainer class object. We provide maximum number of epochs and checkpoint callbacks for retrieving training and validation losses and IoUs for each epoch. The next section discuss more about these metrics.

H. Results

The training loss, validation loss, training IOU and validation IOU(Intersection over union) curves have been plotted for both the models. We also tested the trained model on the images from CityScapes Test dataset and images that we have collected. The results are shown below.

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

432

H.1. UNet and DeepLabV3

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

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

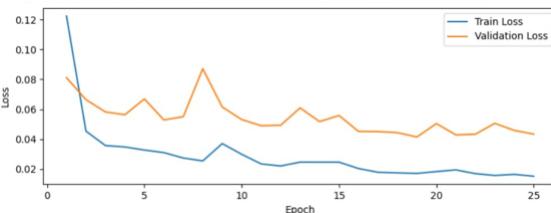


Figure 11. UNet: Training and Validation Losses

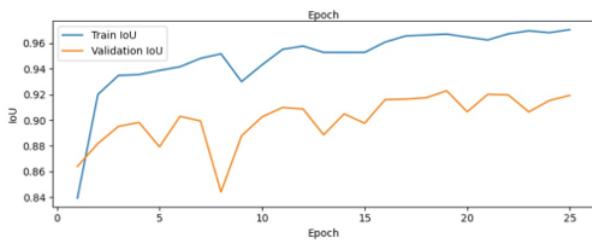


Figure 12. UNet: Training and Validation IOU

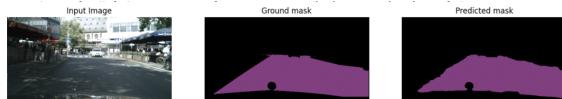


Figure 13. UNet: CityScape Test image1

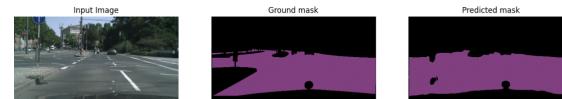


Figure 14. UNet: CityScape Test image2

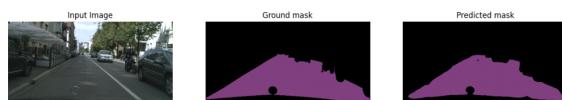


Figure 15. UNet: CityScape Test image3

478

479

H.2. Comparison

480

481

482

483

484

485

As can be seen below in the training and validation loss curves plotted against epochs, the training loss curves turned out to be relatively lesser than that of validation loss, in both models. The curves recovered until the 23rd epoch and started getting higher after that. This is how we fixed the epoch value to 23, to avoid over-fitting. The training

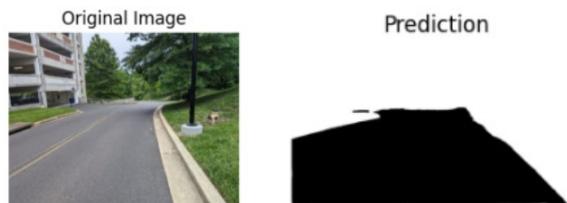


Figure 16. UNet: Our Test image1



Figure 17. UNet: Our Test image2



Figure 18. UNet: Our Test image3

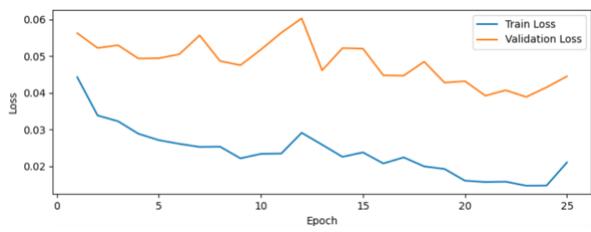


Figure 19. DeepLabV3: Training and Validation Losses

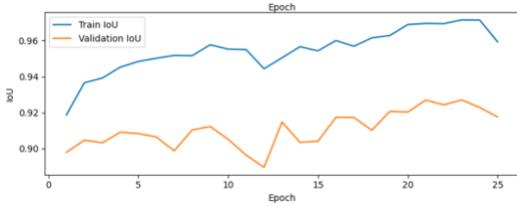


Figure 20. DeepLabV3: Training and Validation IOU

loss for DeepLabV3 is observed to be relatively lesser (Figures 11 and 19). We also observed that the segmentation outputs from the DeepLabV3 model turned out more efficient for our test data. IOU is calculated by measuring the overlap between the predicted bounding box or segmenta-

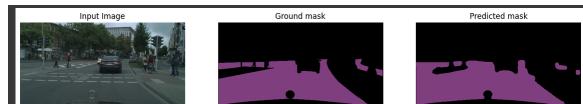


Figure 21. DeepLabV3: Cityscape Test image1

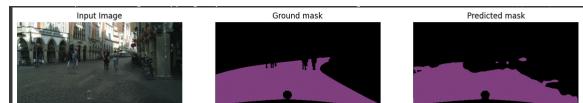


Figure 22. DeepLabV3: Cityscape Test image2

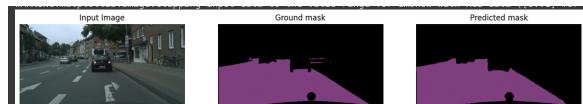


Figure 23. DeepLabV3: Cityscape Test image3

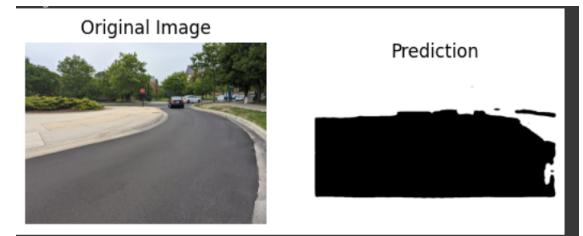


Figure 24. DeepLabV3: Our Test image1

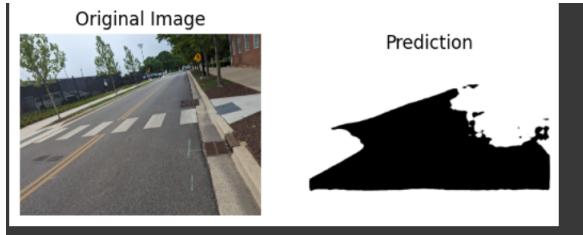


Figure 25. DeepLabV3: Our Test image2

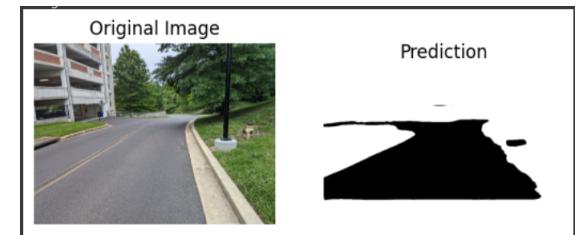


Figure 26. DeepLabV3: Our Test image3

tion mask and the ground truth bounding box or mask. It is computed as the ratio of the area of intersection between the predicted and ground truth regions to the area of their union. The most ideal case is when the IOU is 1. As we can

see in Figure 12 and Figure 20, both the models have their IOU curves almost reaching 1 near the epoch 25. For some of our test cases, both the models are detecting sidewalk as a road which is incorrect. There is a scope of improvement to train the models on custom dataset.

References

- [1] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The Cityscapes Dataset for Semantic Urban Scene Understanding,” in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
- [2] He, Kaiming, et al. ”Deep residual learning for image recognition.” Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [3] Ronneberger, O., Fischer, P., Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image Segmentation. In: Navab, N., Hornegger, J., Wells, W., Frangi, A. (eds) Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015. MICCAI 2015. Lecture Notes in Computer Science(), vol 9351. Springer, Cham. https://doi.org/10.1007/978-3-319-24574-4_28
- [4] Torchvision Semantic Segmentation – PyTorch for Beginners <https://learnopencv.com/pytorch-for-beginners-semantic-segmentation-using-torchvision/>
- [5] HybridNets: End-to-End Perception Network - <https://arxiv.org/abs/2203.09035>
- [6] MIT Scene Parsing Benchmark <http://sceneparsing.csail.mit.edu/>

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