

Robust Lane Detection - Final Report

Charles Downs¹, Arbër Demi¹, Matej Havelka¹ and Georgi Trevnenski¹

¹TU Delft

{C.A.Downs, A.Demi, M.Havelka, G.D.Trevnenski-1}@student.tudelft.nl

1 [Check out our code on our GitHub.](#)

2 1 Introduction

3 Lane detection has, for a long time, been an important topic
4 in the field of computer vision, and has emerged as an in-
5 dispensable tool in modern transportation systems (handling
6 tasks such as lane assistance, auto-pilot, and lane-keeping/
7 assistance). These advances largely hinge on the ability of
8 computer vision systems to accurately and reliably interpret
9 road markings and traffic scenarios, making lane detection a
10 quintessential component in the modernization of transport
11 systems. Progress in these areas of computer vision does
12 not only improve road safety but also traffic fluidity. This
13 progress, however, has proven to be challenging on a variety
14 of levels. While the task of lane recognition and detection is
15 trivial for humans, computer vision systems struggle with as-
16 pects such as environmental changes, occlusions, varying in-
17 frastructure, or faded lanes. Traditional methods, such as im-
18 age processing-based pipelines—which we will further dis-
19 cuss and evaluate in section 2—have proven to be efficient
20 under certain circumstances, but they largely struggle with
21 the previously mentioned challenges. Lane detection systems
22 were improved with the advent of the convolutional neural
23 network, but still struggle under certain circumstances, such
24 as occlusions.

25 The development of self-attention, pioneered by the Trans-
26 former model [Vaswani *et al.*, 2023], has burgeoned recently
27 as a relatively new deep learning architecture for sequence-
28 based tasks. Clearly, there is a sequential nature to lane de-
29 tection tasks, if one assumes a real-time classification task
30 has as input a stream of frames. (Self)-attention mechanisms
31 have proven to be amongst some of the most powerful archi-
32 tectures with sequence-to-sequence tasks, such as natural
33 language processing, translation tasks, and video processing.
34 By treating a stack of images as a sequence, one can tap into
35 the spatiotemporal relationships that often provide valuable
36 insights for lane detection, insights that traditional methods
37 might overlook. As a result, there has been a (small, yet co-
38 ordinated) effort to produce both the first-ever frame-based
39 traffic scene datasets, as well as the first set of deep learn-
40 ing architectures that exploit the inherent spatiotemporal na-
41 ture of this vision task. Of particular interest is the CULane
42 dataset [Xingang Pan and Tang, 2018], a collection of over
43 55 hours of traffic scenes, with a total of 133,235 frames. For

each frame, the authors provide an annotated ground truth. 44
This is in contrast to other datasets, such as TUSimple [Yoo 45
et al., 2020], which also provide a collection of traffic scene 46
footage, however, the authors only provide one annotated 47
ground truth per clip, making it less performant for a truly 48
recurrent architecture. 49

The most recent and relevant work that attempts lane de- 50
tection with a Transformer-based architecture, published in 51
2023, is the O2SFormer [Zhou and Zhou, 2023]. The au- 52
thors use Transformers for a hybrid approach that performs 53
both one-to-one assignments (one detection is assigned to a 54
single ground truth) as well as one-to-several (a single detec- 55
tion can be matched to multiple ground truths). By adopting 56
a one-to-several strategy, the model can more flexibly and ac- 57
curately associate detections with their corresponding ground 58
truths. The authors argue this has the dual advantage of re- 59
solving ambiguity during the training phase and ensuring that 60
each detection is accurately labeled, thereby potentially im- 61
proving the model during real-world deployment. Another 62
significant aspect of the O2SFormer model is the introduc- 63
tion of dynamic adjustments, such as the layer-wise soft la- 64
bel and the dynamic anchor-based positional query. These 65
innovations address the inherent limitations of DETR’s tradi- 66
tional approach, especially in terms of optimizing one-to-one 67
assignment and providing a robust positional context for de- 68
tections, respectively. 69

However, while this model makes a significant im- 70
provement over previous state-of-the-art models (77.83% 71
F1 score on CULane, outperforming existing Transformer- 72
based models), it does have limitations. Specifically, the 73
model does not make full use of the sequential nature inher- 74
ent in video frames or consecutive images (namely, a several- 75
to-one approach). Exploiting the spatiotemporal relationship 76
between sequential frames can potentially provide additional 77
context and information that can enhance detection accuracy, 78
especially in dynamic environments, where road scenarios 79
change rapidly. We, therefore, posit that these Transformer- 80
based approaches can be further improved upon if they adopt 81
a several-to-one approach, where the spatiotemporal nature 82
of the frames is truly exploited. Throughout this paper, we 83
experiment with how Transformers can be used in a several- 84
to-one fashion using the CULane continuous-frame dataset. 85
We investigate multiple stages within the machine learning 86
pipeline. This encompasses enhancements in data prepro- 87

88 cessing (optimizing the data format for better model inter- 142
89 pretability), a diverse selection of backbones for sophisticated 143
90 feature extraction, and the adaptability of Transformers in 144
91 scenarios involving continuous data streams. 145

92 2 Literature Review 146

93 As with many computer vision tasks, lane detection origi- 147
94 nates from traditional image and signal processing pipelines 148
95 [Tang et al., 2021]. Traditional methods would employ a vari- 149
96 ety of normalization/pre-processing to deal with noise and 150
97 other artifacts (such as the usage of mean, median [Wang et 151
98 al., 2010], Gaussian [Hsiao et al., 2009] or FIR filters [Wang 152
99 et al., 2011]) followed by feature-modeling methods such as 153
100 edge detection filters (Sobel and Canny [Tang et al., 2021]) 154
101 as well as more advanced transform-based feature extrac- 155
102 tors, such as the Hough transform [Tang et al., 2021]. Since 156
103 the advent of the convolutional neural network, the epony- 157
104 mous AlexNet, and the explosion in the amount of research 158
105 in deep learning, these novel data-driven approaches have 159
106 largely outperformed traditional methods, as well as elimi- 160
107 nated the need for arduous pre-processing steps such as filter- 161
108 ing, instead performing end-to-end learning, where all under- 162
109 lying noise distributions are being learned and dealt with by 163
110 these deep learning-based methods. More specifically, deep 164
111 learning-based lane detection models can be classified into 165
112 the following categories: segmentation-based, classification- 166
113 based, and regression box-based. In most cases, the latter is 167
114 the most commonly employed method, where the objective is 168
115 to parametrize a shape or region by some parameters (e.g., a 169
116 line can be parametrized by an angle and a magnitude, as a 170
117 vector). This approach is the least computationally expensive 171
118 since it does not require performing per-pixel labels. 172

119 Before the progress in sequence-based models, the corner- 173
120 stone operator was the convolution, used inside CNNs in or- 174
121 der to progressively learn richer features, allowing for tasks 175
122 such as classification. However, a lot of the early break- 176
123 throughs in lane detection were achieved using CNNs. One 177
124 of the early and most widely-cited papers that employ a CNN 178
125 is Line-CNN [Li et al., 2019]. Based on the Faster R-CNN 179
126 architecture, Line-CNN replaces the *Region Proposal Unit* 180
127 (RPU), with a *Line Proposal Unit* (LPU), effectively pro- 181
128 ducing a line rather than a region. Line-CNN predicts k line 182
129 shapes at each location for the boundary (bottom, left, right), 183
130 where each line has 2 scores for the confidence of being a real 184
131 traffic lane. The LPU itself is implemented as a Fully Convo- 185
132 lutional Network (FCN), where the last convolutional feature 186
133 map is used in order to extract k line proposals. A 1024- 187
134 dimensional feature vector is extracted from each boundary 188
135 position in the feature map. This feature map is then fed into 189
136 two sibling classification and regression networks, which in 190
137 turn output a set of k lines, each with output scores (the scores 191
138 mentioned above) and then a set of $k \times (S + 1)$ coordinates 192
139 for each line $k_i \in k$. This FCN effectively acts as a fea- 193
140 ture extractor backbone, the authors here use a ResNet, with 194
141 122 layers¹ as their backbone of choice for the feature extrac- 195

¹The authors argue that with 122 layers, the receptive field of the last convolutional feature map is large enough to cover the whole image, which matches traffic lines

142 tion. In terms of model architecture, Line-CNN is archetypal: 143
144 it combines a feature extractor, followed by the main model 145
146 which performs the detection task. In terms of performance, 147
148 the authors benchmark their model on the MIKKI dataset, and 149
149 the TuSimple dataset. The authors evaluate their method 150
151 against 3 other methods, of which only 2 employ deep learn- 152
152 ing. As this paper was one of the first works to employ deep 153
153 learning for lane detection, it naturally achieved state-of-the- 154
154 art results, having better accuracy, precision, and recall than 155
155 the other methods, the details of which we spare from this 156
156 report. 157

158 With this in mind, we can now introduce models that 159
159 employ a recurrent/sequential structure in their architecture. 160
160 With the introduction of recurrent sequence-based architec- 161
161 tures, we can truly exploit the inherent nature of traffic scenes. 162
162 Of particular interest is the work from [Zou et al., 2019], 163
163 with their introduction of ConvLSTM mechanisms, which 164
164 can capture temporal information—making it fitting for a lane 165
165 detection task. The overall architecture of this model fol- 166
166 lows an encoder-decoder structure, with ConvLSTM in the 167
167 bottleneck. The encoder is used as a feature extractor, simi- 168
168 lar to in Line-CNN, which then feeds the feature maps into 169
169 the RNN blocks. The encoder-decoder network then mod- 170
170 els the lane detection as a semantic segmentation task (as is 171
171 common with encoder-decoder type networks). The decoder 172
172 then employs deconvolution and upsampling to highlight the 173
173 targets and spatially reconstruct them. Both the encoder and 174
174 decoder are Fully Convolutional (FC). In terms of training, 175
175 the authors train the network using the pre-trained weights 176
176 from SegNet and U-Net trained on ImageNet (a relevant as- 177
177 pect of this paper, as we will discuss pre-trained networks for 178
178 lane detection in the methodology section). The authors em- 179
179 ploy a cross-entropy loss to treat the task as a discriminative 180
180 segmentation task. The authors use the TuSimple dataset, 181
181 as well as their own dataset, one noteworthy aspect (particu- 182
182 larly in regard to our own model) is that both these datasets 183
183 only contain 1 ground truth image per video, meaning that 184
184 the overall accuracy and loss will not always be entirely ac- 185
185 curate, this is in contrast to the CuLane dataset, which has a 186
186 labeled ground truth for each frame in a clip. The authors 187
187 then evaluate their results using precision, recall, and F1- 188
188 measure. The authors benchmark their model against a vari- 189
189 ety of encoder-decoder modules, such as SegNet and UNet, 190
190 as well as publish the results of their model compared to other 191
191 benchmarks on the TuSimple dataset, where their model is 192
192 second in terms of accuracy (as of March 14th, 2018). Fol- 193
193 lowing a similar architecture as above, the authors of [Dong 194
194 et al., 2023] propose a comparable encoder-decoder model, 195
195 using what they call a “spatiotemporal RNN”, or ST-RNN in 196
196 order to learn spatiotemporal (ST) dependencies. The authors 197
197 argue that most existing methods do not take full advantage 198
198 of the essential properties of the lane being long continuous 199
199 solid dashed lines, nor do they make use of the ST informa- 200
200 tion correlation and dependencies in the continuous driving 201
201 frames. The authors treat the detection task as a segmen- 202
202 tation task, in contrast to [Zhou and Zhou, 2023], however, 203
203 the authors use a novel hybrid sequence-to-one deep learning 204
204 architecture. Additionally, the authors propose an encoder 205
205 CNN with SCNN [Pan et al., 2018] layers, and then a decoder 206

201 CNN with FC layers. Here, SCNN refers to what is known
 202 as a “spatial CNN”—an approach to capture spatial relation-
 203 ships using convolution. To better exploit the spatiotempo-
 204 ral relationships, the authors introduce the aforementioned
 205 SCNN layer in the encoder to extract and make use of these
 206 relations in a given image frame. The features are then fed
 207 into the ST-RNN module, the authors employ both ConvL-
 208 STM and ConvGRU to model the time-sequence aspect, simi-
 209 larly to [Zou *et al.*, 2019]. The novelty here is the usage of
 210 the SCNN module, which performs a more robust feature ex-
 211 traction by propagating spatial information via slide-by-slide
 212 message passing (UP, DOWN, LEFT, RIGHT). The authors
 213 benchmark their method against a variety of other mod-
 214 els, including SCNN, SegNet, UNet, and LaneNet. The au-
 215 thors try a variety of backbones as well as RNN architectures,
 216 resulting in 13 different models, ranging from fairly shal-
 217 low and lightweight to deeper SegNets. Overall, the authors
 218 achieve the best performance compared to the other models
 219 (using common scores such as F1, precision, and recall).

220 Having discussed some of the sequence-based models, we
 221 can now introduce the concept of attention mechanisms. As
 222 discussed in section 1, attention mechanisms have become
 223 a very popular choice for sequence or time-series data, with
 224 their strong ability to identify relationships in data, making it
 225 a good fit for lane detection. This was explored by [Tabelini *et*
 226 *al.*, 2021] in a relatively recent paper. The overall architecture
 227 of this approach also employs a CNN for feature extraction,
 228 in this case, the authors use a ResNet. Here the authors op-
 229 erate on what they call *anchor points*, which are obtained by
 230 pre-processing the lane annotations (the collection of points
 231 $\{(x_i, y_i)\}_{i=0}^N$) such that they obtain anchors defined by i) an
 232 origin $O = (x_{orig}, y_{orig})$, with y_{orig} located on the border
 233 of the y -axis, and then ii) a direction θ . In this paper, this
 234 comprises a total of 2782 anchors. The feature extraction
 235 operates directly on this data, where the feature map corre-
 236 sponds to extracted features from each anchor. These fea-
 237 tures are extracted through a process the authors call *anchor-*
 238 *based feature pooling*, where a feature vector for a specific
 239 anchor is pooled from a feature map. Here, an anchor de-
 240 fines the points of the feature map that will be used for the
 241 respective proposals. We note that since these anchors are
 242 modeled as lines, the interest points for a given vector are
 243 those that intercept the anchors *virtual line*. This is similar to
 244 Line-CNN, however, the authors do not fix the y coordinate
 245 to be clipped to the bottom border, yielding more profound
 246 features. Due to the intricacy of the next part, we propose
 247 to introduce some notation: the feature vector extracted from
 248 the feature map \mathbf{F} will be named \mathbf{a}_i . Due to the nature of
 249 the feature pooling (and the inherent nature of CNNs) the in-
 250 formation carried by the feature vector is mostly local. To
 251 overcome this, the authors introduce attention mechanisms
 252 to produce global features, \mathbf{a}^{glob} , which are aggregated with
 253 the local features. This attention mechanism is composed of
 254 a fully-connected layer, L_{att} , which processes the local fea-
 255 tures \mathbf{a}_i , and outputs a weight (via a softmax operator) for
 256 how much attention each local feature should get. This is
 257 denoted as $w_{i,j}$ to represent the attention weight from posi-
 258 tion i to position j . \mathbf{a}^{glob} is then aggregated with the local

259 features by a weighted summation, as is the case with classi-
 260 cal attention mechanisms. A lane proposal is then predicted
 261 for each anchor, which is then fed into two FC layers; one
 262 for classification (probability of proposal i being a lane), and
 263 another for regression, which parametrizes the line. The au-
 264 thors then employ non-maximum impression (NNMS) in or-
 265 der to reduce the number of false positives. In terms of ex-
 266 periments, the authors train and test their model on 2 notable
 267 datasets: CuLane and TuSimple, both of which are rele-
 268 vant to this research project. **TuSimple**: In terms of accu-
 269 racy, the author’s method is on par with other state-of-the-art
 270 methods (including SCNN, Line-CNN, and PolyLaneNet).
 271 The authors, however, note that these datasets mostly consist
 272 of simple cases, where the metrics are rather permissive.
 273 The authors, however, note the speed difference with their
 274 own method is significantly higher than other methods. **Cu-**
 275 **Lane**: Lastly, when benchmarking against CuLane, the au-
 276 thors note that their model achieves the highest F1 amongst
 277 all compared methods while keeping a relatively high FPS.
 278 The model also performs well considering that this dataset
 279 contains more challenging and complex scenes.

280 Finally, we propose to conclude this review with some of
 281 the latest papers that focus on lane detection using Transfor-
 282 mers. More specifically, [Zhou and Zhou, 2023] proposed a
 283 One-to-Several architecture in 2023. As their premise, the
 284 authors note that although some attempts have been made
 285 to employ Transformers (more specifically DETR, or DEtec-
 286 tion TRansformer [Carion *et al.*, 2020]) for lane detection,
 287 two main issues that limit their performance are i) the po-
 288 sitional query in DETR lacks a clear focus area, making it
 289 difficult to detect lanes, and, ii) the one-to-one label assign-
 290 ment causes low training efficiency due to label semantic con-
 291 flicts. To overcome each of these concerns, the authors pro-
 292 pose the following: i) the authors design a dynamic anchor-
 293 based positional query, where lane anchors are encoded as
 294 a positional query. This allows the exploration of explicit
 295 positional prior to the positional query. Next, they intro-
 296 duce One-to-Several label assignments to solve label seman-
 297 tic conflicts, this works by using one-to-many assignments for
 298 the first $N - 1$ layers of the decoder, and then one-to-many
 299 for the last layer. Lane anchors are represented in a simi-
 300 lar fashion to [Li *et al.*, 2019], where x -coordinates have a
 301 direct mapping to the y -coordinates, based on a uniform sam-
 302 pling along the edge of the image, and the network output is
 303 identical too (length of lane anchors, start and end point, an-
 304 gle, and probability of existing). The setup of the model is
 305 archetypal, as it follows the widely-cited DETR model, using
 306 a pre-trained ResNet as a backbone. The authors evaluate the
 307 model on the CuLane dataset, using the F1 score as a preci-
 308 sion metric. Their model is then evaluated against many state-
 309 of-the-art models, including SCNN and LanteATT, discussed
 310 above. The authors benchmark their model with ResNet18,
 311 ResNet34, and ResNet50, with the expectation that the deeper
 312 ResNets achieve higher scores. The authors benchmark their
 313 model in a variety of challenging scenarios that are avail-
 314 able in the CuLane dataset², where even with ResNet18, the

²Normal, Crowded, Dazzle, Shadow, No Line, Arrow, Curve, Night

315 model outperforms almost all other models in all scenarios, 316 with the ResNet50 backbone achieving the absolute highest 317 F1 score across all scenarios, and against the 8 other mod- 318 els (LSTR, SCNN, UFLD, PINet, CurveLane-L, LaneATT, 319 LaneFormer with ResNet18 and ResNet34 backbone). The 320 authors also examine the convergence curves for their model 321 versus vanilla DETR, noting that their model converges much 322 faster than DETR, illustrating the clear improvement over 323 DETR.

324 3 Methodology

325 Having discussed pre-existing research, as well as their meth- 326 ods and results, we would like to discuss the approach and 327 methodology employed in this project. First and foremost, 328 for a better comprehension of how the project evolved based 329 on discoveries and results, we heavily encourage the reader 330 to refer to the methodology section from the mid-term report. 331 For ease of access, it has been added to the appendix, in sec- 332 tion A.

333 The main architectural changes involved switching from an 334 anchor-based approach (analogous to those discussed in sec- 335 tion 2) to a segmentation-based approach and replacing the 336 ResNet feature extractor with a model that—following our 337 investigation—seems more fitting. The Transformer-based 338 approach (namely, exploiting self-attention) remains present 339 in our approach. We propose to outline each of these 2 340 changes in the subsections below, so as to give sufficient de- 341 tail into our feature extractor, followed by our main model 342 and architecture. These 2 aspects naturally follow each other, 343 but before we start on the changes made, we will dive into the 344 dataset used for this work.

345 3.1 Dataset

346 The lane detection dataset that we used for our experimen- 347 tation is CULane [Xingang Pan and Tang, 2018]. As men- 348 tioned before, the dataset is comprised of 133235 frames, 349 adding up to 55 hours of video. The frames are of 1640×590 350 resolution. The dataset is split into 88880 frames for the train- 351 ing set, 9675 for the validation and 34680 for the test set. The 352 test set has its own separation into different categories, each 353 with a separate challenging scenario and making up a differ- 354 ent percentage of the test set:

355 Normal, Crowded, Night, No line, Shadow, Arrow, Dazzle 356 light, Curve, Crossroad.

357 Examples of each scenario can be seen in the appendix 8. 358 The lanes for each frame are annotated with cubic splines, 359 which are just piecewise cubic functions that interpolate a set 360 of datapoints.

361 A worthwhile mention is also the fact the annotators chose 362 to annotate lanes in some frames where lanes are not visible, 363 but there is enough context from the video to know that they 364 are there. At the same time, some lanes that are visible are 365 not annotated, such as lanes behind barriers, where entry in 366 the visible area with a vehicle would not be expected.

367 The main focus of the dataset is four main lane markings 368 which are most used in real-world scenarios, therefore other 369 lane markings are not annotated.

370 However, this main focus does not seem to always be con- 371 sistent. There are cases where annotations for lanes are made

372 but there is not really any visual context other than general 373 driving experience that shows that lanes should exist in that 374 particular location. This can be seen in the appendix 10. 375 These cases can be damaging to the training process depend- 376 ing on the pipeline used.

377 At the same time, there are also examples of lanes being 378 fully visible, however they are annotated only after the car 379 starts moving from a stop, and not before. This makes a lot 380 of frames with lanes present appear as if there are no lanes in 381 the ground truth. Two such cases are presented in appendix 382 9.

383 Some practical information about the dataset that is not 384 made immediately explicit by the authors is that the videos 385 provided are not all at the same sampling rate, some being 386 sampled every 30 frames, and some every 90 frames, which 387 might have effects on performance depending on the pipeline 388 being used.

389 3.2 Feature Extraction

390 Most research on lane detection, including the studies cited in 391 our literature review, utilizes a pre-trained network in order to 392 perform low-level feature extraction—commonly a ResNet or 393 U-Net. Typically, these feature extractors (herein referred to 394 as backbones) are pre-trained on large datasets such as Im- 395 ageNet. The ImageNet dataset consists of a large number 396 of images and has been instrumental in advancing computer 397 vision. The wide variety of images featured in this dataset 398 makes it an attractive feature extractor for a variety of down- 399 stream tasks, plainly due to the fact that it has the ability to 400 capture the most salient features in many downstream vision 401 tasks. In the case of lane detection, however, it would not 402 be unreasonable to doubt the saliency and descriptiveness of 403 these features. Lanes are typically described as long white 404 lines, any reasonable feature extractor for downstream tasks 405 should be able to capture these markings.

406 In an effort to investigate the saliency of these features, 407 we set up a variety of experiments in order to help us eval- 408 uate feature extractors. Given that the ResNet architectures 409 are some of the most widely-used backbones for lane detec- 410 tion, we trained a ResNet34 in an auto-encoder setup. We im- 411 port a ResNet34 with ImageNet default weights, and omit the 412 last fully connected layer, as we are not interested in a clas- 413 sification task. The modified ResNet34 acts as the encoder, 414 outputting a feature vector of size $[1, 512, 1, 1]$, cor- 415 responding to the $[N, C, H, W]$ convention. This effec- 416 tively gives us a 512-sized vector as per the ResNet34 archi- 417 tecture specification, representing the latent space of the input 418 image. The encoder’s task ends with the generation of this 419 latent representation. The subsequent challenge is for the de- 420 coder to reconstruct the original image from this compressed 421 form. The success of the decoder in recreating the lanes from 422 the latent representation would demonstrate the sufficiency of 423 the extracted features for the image reconstruction task. This 424 feature vector is then fed into a decoder layer, effectively con- 425 stituted of a succession of 2D transposed convolutions, and 426 ReLU activations in order to re-create the image. This model 427 was then trained on a subset of the CULane dataset over the 428 course of 12 hours, for a total of 100 epochs. These results 429 are illustrated in the appendix, figure 11. These images are

430 both results obtained during the later epochs in the training
431 stage. The lower picture illustrates the best result obtained,
432 where a car can be (somewhat) identified (see evaluation and
433 appendix). In the other case, the most salient features cap-
434 tured are the blue sky and the grey road. Little to no lane can
435 be distinguished, and it is noteworthy that these images are
436 hand-picked from the training set, most likely indicating the
437 validation or real-world reconstructions would not necessar-
438 ily be better. This can be seen as an indication that ResNet
439 (in our case, ResNet34) is perhaps not the most adequate fea-
440 ture extractor when used pre-trained on ImageNet. Consid-
441 ering the results obtained here are after specific training on a
442 lane detection dataset, it is safe to assume that these types of
443 backbones are sub-optimal for this task. With this in mind,
444 we opted to investigate alternative models in order to create a
445 better feature extractor. We followed the same auto-encoder
446 approach as mentioned above in order to validate the effi-
447 cacy of the feature vector, where a decoder would be tasked
448 with re-creating the images from the feature vector. Due to
449 the inherent sequential nature of these datasets, self-attention
450 mechanisms may be a better fit, as one would expect if a lane
451 is there, it would be present across multiple frames, despite
452 the environment changing. With this in mind, we decided to
453 train a Vision Transformer (ViT) encoder, once again with
454 the MLP head removed. This approach proved sub-optimal,
455 as the sequential nature of the frames being fed to the ViT
456 would cause the model to simply learn to output a clone of
457 the ground truth it received after multiple frames. A more
458 robust approach would use a masked auto-encoder (MAE)
459 [He *et al.*, 2022], such that the feature extractor is forced
460 to learn a representative latent feature space. The approach
461 here is to load the MAE’s pre-trained weights from ImageNet
462 into an autoencoder ViT, which is effectively an autoencoder
463 that employs the ViT architecture for its encoder and decoder.
464 We then train the decoder to reconstruct the image from the
465 latent space, with frozen weights on the encoder—as un-frozen
466 weights were not possible given time and hardware constraints
467 of this project. We use a masking ratio of 0, however, the
468 weights of the imported encoder were trained with a 75%
469 masking ratio, as per the original MAE paper, meaning the
470 added robustness from the masking of the reconstruction is
471 still present in this downstream task, yet we do not withhold
472 any information from our model itself, as missing patches can
473 easily break the continuousness of lanes, thereby perturbing
474 learnable features. A sigmoidal activation is then placed on
475 the output of the decoder in order to produce pixel values in
476 the $[0, 1]$ range. The results are shown in figure 13 and 15
477 in the appendix, with pairs representing ground truth and out-
478 put, respectively. It is worthwhile mentioning here that the
479 learned features here are not from CULane, in the same way
480 as we did with ResNet34, but rather from the pre-training
481 performed on ImageNet. Better results could have been
482 achieved if the encoder had been trained on the CULane dataset
483 with a small masking ratio (e.g. 25%), and no masking for
484 the fine-tuning phase. However, the computational and time
485 requirements would not have fit in the scale of this project,
486 given the time needed to train full MAEs. We believe that
487 this approach may be more successful for lane detection mod-
488 els that rely on feature extractor networks, by training a feature

489 tractor on the CULane dataset, the model may learn more
490 representative features for the task at hand, resulting in bet-
491 ter performance on downstream tasks. In the state-of-the-art
492 papers reviewed, none had performed pre-training on traffic
493 scene datasets. Results are further discussed in the evaluation
494 section.

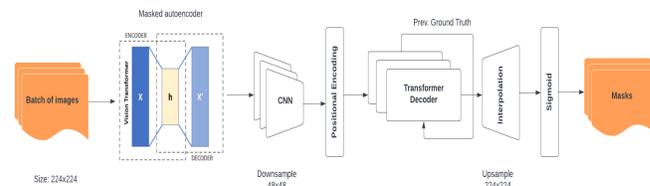
3.3 Main Pipeline and Model

Data processing

495 Before we could start training our model on a segmentation
496 task, we had to preprocess the ground truth CULane provides.
497 The dataset’s annotations consist of a list of lanes per image,
498 which are a collection of (x, y) coordinates of points form-
499 ing a lane. To turn such representation into a segmentation
500 mask, we start with a blank image of the same size as the in-
501 put and draw short straight lines between every two consecu-
502 tive points in a lane for all present lanes. We set a thickness
503 of 5 pixels on these lines which proved to be working well in
504 our experiments. This image we turn into a binary segmentation
505 mask.

506 Furthermore, we performed data augmentation following
507 the example of the O2SFormer. We make use of random hori-
508 zontal flips and random rotations with an angle of rotation
509 varying in a range of 30 degrees. Additionally, we crop the
510 images at random locations and resize them to 224 by 224
511 pixels before they enter our model’s pipeline. A step-by-step
512 visualization of the augmentation process can be seen in Fig-
513 ure 20 in the appendix.

Architecture



514 Figure 1: Pipeline of our model (Partially adapted from [Hinton and
515 Salakhutdinov, 2006]). For more detailed explanation of the pre-
516 trained MAE used please refer to [He *et al.*, 2022]

517 The architecture of our lane detection algorithm consists of
518 two major components - a Masked Autoencoder (MAE) ([He
519 *et al.*, 2022]) and a decoder layer. As Figure 1 displays, be-
520 tween these two components, we utilize a standard convolu-
521 tional neural network to downsample the images so they can
522 fit on the GPU we used to train. The CNN consists of three
523 convolutional layers with average pooling layers in between.
524 The downsampled images are then positionally encoded be-
525 fore they are fed to the transformer decoder. The positional
526 encoding used is a traditional one with sine and cosine func-
527 tions, suggested by the original Transformer model ([Vaswani
528 *et al.*, 2023]).

529 The MAE is used in our architecture as a feature extractor
530 and for this purpose, we load the entire model, pre-trained on
531 the ImageNet dataset. This autoencoder works by dividing
532 the input image into patches (squares of a specific size) and

533 a percentage of those are masked - set to a pixel value of 1.
534 The unmasked patches are run through an encoder and the
535 encoded patches are assembled together with mask tokens in
536 a feature map with the size of the original image. A decoder
537 then tries to reconstruct the original image given the encoded
538 features. Important to note here is that masking is used only
539 in the pertaining of the MAE and we do not mask the images
540 in our pipeline. The way we incorporate the already trained
541 MAE in our architecture is by freezing its encoder's weights
542 and training only its decoder in our backward pass. There-
543 fore, the decoder is expected to learn to output a feature map
544 optimized for the segmentation task we have.

545 The decoder we use after the MAE is a standard PyTorch
546 implementation without any LayerNorms of a transformer de-
547 coder with one decoder layer. The encoded feature maps cor-
548 responding to the images in the batch are fed into it one by
549 one. The ground truth of the previous image is passed as a
550 target sequence to the decoder during training, except for the
551 first image where we input a tensor of zeroes. This is the
552 conventional way of training transformers for NLP tasks and
553 even though not recommended in segmentation tasks. In the
554 end we decided to follow this, as to help the network to learn
555 from previous frames. We also tested an approach where we
556 stop feeding it the ground truth once the validation loss stops
557 falling, which gives promising results. During evaluation, of
558 course, the output of the decoder from the previous iteration
559 in the batch is used instead of the ground truth.

560 Finally, the output of the transformer decoder is upsampled
561 to the original size of the images through bilinear interpola-
562 tion and a sigmoid function is applied to it. The end result
563 per image is a segmentation mask with a probability for ev-
564 ery pixel whether it belongs to a lane or not.

565 3.4 Loss

566 Throughout the time we have worked on our pipeline, we
567 have experimented with many different losses in attempts to
568 rectify the issues we were faced with. We will mention them
569 here, however, most will be brief as they were used at a time
570 when our model was not fully functional for reasons unrelated
571 to the loss.

572 IoU Loss

573 The main loss we used during training for our final version
574 of the pipeline is a self-made IoU loss. It is quite simple,
575 however, it has proven to be effective. The loss is formed by
576 components in the numerator and denominator, *up* and *down*,
577 where

$$up = \sum_{dim=1,2} prediction * target$$

578 with prediction being the result of the model and target being
579 the ground truth. Whereas

$$down = \sum_{dim=1,2} 0.5 * prediction + 0.5 * target$$

580 This makes the IoU differentiable, while still keeping the
581 overall expected shape of the function. The main reason for
582 picking this loss is that it is ensured that the global minimum
583 is only at the point where the prediction perfectly matches the

ground truth.
The end result is

$$IOU_{loss} = 1 - \frac{1}{N} * \sum \frac{up + \epsilon}{down + \epsilon}$$

584 where $\epsilon = 0.000001$ to address some frames having no lanes,
585 which would lead to division by zero.

588 Cross Entropy Loss

589 Another loss that we worked a lot with was the PyTorch im-
590 plementation of the CrossEntropyLoss, which can be found
591 in appendix D.1. We tried a weighted approach as we had
592 quite an unbalanced training set, with many cases being an-
593 notated as lanes, and few without. Although it was promising,
594 most of the experimentation was during the faulty stage of our
595 pipeline, so we cannot make conclusions on its performance
596 for our case.

597 We also tried the Binary CE loss from PyTorch which also
598 showed promising results, however, it also was in the early
599 stages of our pipeline.

600 4 Evaluation

601 In order to evaluate results, we propose the same structure
602 as the methodology; where the section is split into the fea-
603 ture extractor evaluation, and then the main model backbone
604 itself. One important aspect to note is that, given time and
605 hardware constraints, the authors were unable to train the
606 model until convergence. All training was done on the au-
607 thor's own personal hardware. Given the size of the dataset,
608 and the time required for convergence, a realistic comparison
609 and a fully-converged model are not realistic. We therefore
610 propose to evaluate the results obtained *so far* in an *a pri-*
611 *ori* manner. Training is particularly cumbersome due to the
612 MAE.

613 4.1 Feature Extractor

614 An appropriate evaluation of the MAE ViT feature extractor
615 would be to benchmark it against some of the other state-of-
616 the-art methods discussed in the literature review. Unfortu-
617 nately, due to time and hardware constraints, this was not pos-
618 sible. All training was performed on the author's own, per-
619 sonal hardware. On the whole, the auto-encoder setup used in
620 order to train the MAE ViT seemed to show promising results
621 on image re-creation, with lanes being identifiable in almost
622 all cases. The only comparison we can do is to visually in-
623 spect how ResNet34 performed in a similar setup, obviously,
624 we have not tested the feature extractor on downstream tasks.
625 Visual comparisons of reconstructed images can be seen in
626 figures 11 and 13, where the ResNet34 output seems to only
627 be able to re-create road and sky, whereas the proposed model
628 seems to re-create the entire image.

629 4.2 Main model

630 To conclude, we will discuss the results obtained on the main
631 model and architecture. As discussed above, the feature ex-
632 tractor backbone is using the ImageNet pre-trained weights.
633 Once again, due to time and hardware constraints, we were
634 not able to train the model long enough to get it to converge,
635 as can be seen in the loss graphs below.

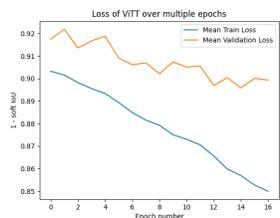


Figure 2: Loss over 16 epochs, training enabled (ground truth being fed back into decoder, see 1.

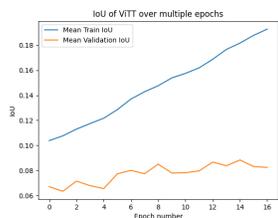


Figure 3: IoU over 16 epochs, training enabled (ground truth being fed back into decoder, see 1.

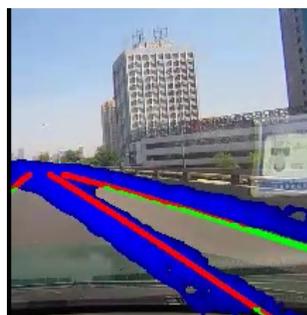


Figure 4: Loss over 16 epochs, without ground truth fed into decoder



Figure 5: IoU over 16 epochs, without ground truth fed into decoder



Figure 6: Loss over 16 epochs, without ground truth fed into decoder



Figure 7: Loss over 16 epochs, without ground truth fed into decoder

636 As illustrated with the IoU and loss on the test set, the model is clearly not converging. We intertwined the running
 637 of training and validation in order to evaluate the model at each epoch. The training time for 16 epochs was roughly
 638 30 hours, a significant amount of time considering that we have been advised that models only start learning after 10
 639 epochs³. On these types of tasks, it is not uncommon to train for over 100 epochs, which is simply not viable given the
 640 aforementioned time and hardware constraints. Below (Figures 4,5,6,7) are some illustrations of the output with training
 641 set to **off** (the ground truth is no longer fed to the forward pass, instead model uses its own predictions) after 4 epochs.
 642 Green is ground truth, red is where IoU matches, and blue is output from the model.

643 Clearly, the model is learning to output lanes, as demonstrated by the array of frames outputted, straight lanes seem to
 644 be somewhat easy to capture, but the model still seems to output lanes where there are none (Figure 6) and struggles with
 645 curved lanes (Figure 7), though the segmentation mask does seem to follow the general direction. This is further illus-
 646 trated by the mean IoU which is continually decreasing across frames and has not yet converged. These results do, however,
 647 suffice to demonstrate this approach and method as a proof of concept for this architecture. We have adopted a mixed
 648 ViT/Transformer Decoder approach, combining a novel type of feature extractor, as well as the Transformer decoder using
 649 previous frames to truly exploit the sequential nature of the problem.

650 Finally, an important limitation that we faced due to the hardware constraints is the high batch size required (32), on
 651 the relatively large images used. As discussed in the methodology section, our architecture had to employ a convolutional
 652 neural network (CNN) in order to downsample the images such that they could fit on our hardware (see appendix
 653 for a full list of hardware specifications and hyperparameters). This has the effect that the up and down-sampling from
 654 48×48 to 224×224 means that a single pixel goes into a roughly 5×5 grid, where lanes have a thickness of 5. This
 655 results in an inherent limitation in how good of an approximation the network can get, but unfortunately was an unavoi-
 656 dable step, considering 1) hardware limitations, 2) the ViT is pre-trained on images of size 224×224 .

657 In Table 1 you can find all the average metrics on the different models we have trained until now. These are ran on

658 ³As noted by this project supervisor

Threshold	IoU	Accuracy	Precision	Recall	F1
0.5	10.4%	84.6%	10.9%	74.1%	9.3%
0.85	11.9%	92.8%	14.7%	45.9%	11.2%

Table 1: Results of the Main model with different thresholds for the statistic computation. With threshold 0.5 we predict a lane if the segmentation mask outputs a value above 0.5 for a pixel.

659 the CULane test set, with the same settings as the model, and feeding of ground truth turned off. We run it on different
 660 thresholds, used to determine whether a lane is predicted or not, based on the results obtained in Figure 18.

661 5 Future work

662 5.1 Content queries

663 Inspired by the concept of content queries described in the O2SFormer paper, we worked in the direction of adding more
 664 positional prior information to the model. Similarly to the approach of O2SFormer, we make use of the DETR model
 665 ([Carion *et al.*, 2020]), which was initially developed to predict sets of bounding boxes. It utilizes a ResNet-50 backbone
 666 trained on ImageNet, a transformer architecture and a shared feed forward network that predicts the class of the detected
 667 object or the lack of such an object.

668 In order to adapt the DETR model to our use case, we drop the feed forward network at the end of its pipeline and use
 669 the feature vectors outputted by the decoder. As we use the loaded pre-trained model we also had to configure the number
 670 of feature vectors.

699 of object queries which are the input to the decoder alongside
700 the encoder output. Each of these object queries specializes in
701 finding information about objects in different parts of the image
702 while going through the decoder layers. Thus, we take the
703 feature vectors from the decoder output and incorporate them
704 with the output of the MAE in our architecture. This setup
705 is different from the one in O2SFormer where the DETR decoder
706 output is fed into the decoder of the O2SFormer architecture
707 as a target sequence. As mentioned before, we use
708 the ground truth of the previous frame as the target sequence
709 and the results from this cannot be outdone by the content
710 queries. Unfortunately, we could not get good results from
711 the addition of the DETR decoder output. We speculate that
712 more fine-tuning and longer training are needed to make the
713 model converge with this setup.

714 5.2 Training optimization

715 Another possible improvement that can be done in the future
716 is an optimization of the training process. Currently, our
717 model starts exhibiting signs of overfitting, i.e. significant
718 difference between the training and validation losses, after a
719 certain point in the training. This problem is overcome
720 after some time and we see the validation loss falling to values
721 close to the training loss. However, an even more optimal
722 training could be performed if we stop feeding the ground
723 truth from the previous frame to the decoder once the overfitting
724 starts occurring and instead use the output of the decoder
725 from the previous iteration. As shown in section 4, once we
726 do not feed the ground truth, the validation loss decreases and
727 continues decreasing in a more stable manner.

728 5.3 Network simplification

729 Lastly, one can simplify our architecture which might lead to
730 quicker and more stable results. Because we did not have the
731 time to run our model on a more powerful machine, we had
732 to include down and up sampling in our architecture, as we
733 could not afford so many parameters. One could possibly
734 remove the down and up sampling which should make it easier
735 for the network to become more precise. Additionally, one
736 could further simplify the MAE decoder to output a smaller
737 dimension, but this would require further work into how to
738 unshuffle the shuffled latent space into a smaller output
739 dimension.

740 5.4 Feature Extractor

741 As mentioned in section 4, we are unable to provide a full
742 evaluation of the feature extractor in downstream lane detection
743 tasks, primarily due to time and hardware constraints.
744 When training the feature extractor on 30% of the data, and
745 with a batch size of 8, preliminary results seem to indicate
746 a better (more rapidly decreasing) loss than the MAE ViT
747 pre-trained on ImageNet. Visual inspections when using the
748 model on the downstream task did not yield satisfactory
749 results. This would ideally be remedied by training the feature
750 extractor on the full dataset, and on a dedicated GPU cluster.
751 If we are to believe that reconstruction is a good indicator of
752 a feature extractor, then this MAE ViT may significantly out-
753 perform other state-of-the-art models using ResNets as back-
754 bones.

6 Ethical and Technical Implications

755 Lastly, we would like to draw attention to the ethical and
756 technical implications that should be considered with such a
757 project. Due to the vast amounts of data required to train con-
758 temporary deep learning models, the quality, as well as quan-
759 tity, of the data needs to be considered. It has been proven
760 that deep learning models suffer from inherent bias present in
761 the training data [Roselli et al., 2019] [Mehrabi et al., 2022],
762 with bias leading to unintended outcomes, such as models
763 that may have a racist or a sexist bias. As outlined in section
764 1, lane detection is inherently difficult due to the lack
765 of consistency in lane markings, lane marking quality, per-
766 country differences in road structure, and environmental con-
767 ditions. Given this complexity, it is imperative to ensure the
768 dataset used for training encompasses a broad range of sce-
769 narios, conditions, and geographical locations (worldwide,
770 ideally). For instance, much of the research in deep learn-
771 ing emanates from developed, wealthier countries, meaning
772 that the model is trained on predominantly well-marked, well-
773 maintained roads. This is exactly the case with the dataset
774 used throughout this project; the CULane dataset consists ex-
775 clusively of clips from the Beijing area, resulting in a two-fold
776 ethical consideration. First, this means that any model trained
777 on this data could be unsafe to operate in regions that have ge-
778 ographical differences and infrastructure differences without
779 extensive testing on a local dataset. Beyond this, this has the
780 effect of making these models only usable in wealthier, de-
781 veloped countries. Under-developed countries, which suffer
782 from worse infrastructure cannot benefit from the develop-
783 ment of these lane detection models if they are exclusively
784 trained on data from predominantly wealthy countries, mak-
785 ing access to this technology unattainable to underdeveloped
786 countries.

787 Privacy aspects can also be a consideration, the function-
788 ality of the model should serve only to predict lane position
789 given a fixed number of frames K . Image frames should not
790 be stored or used beyond the lifecycle needed to predict the
791 lane markings.
792

793 References

- 794 [Carion et al., 2020] Nicolas Carion, Francisco Massa,
795 Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov,
796 and Sergey Zagoruyko. End-to-end object detection with
797 transformers, 2020.
- 798 [Dong et al., 2023] Yongqi Dong, Sandeep Patil, Bart van
799 Arem, and Haneen Farah. A hybrid spatial-temporal deep
800 learning architecture for lane detection. *Computer-Aided
801 Civil and Infrastructure Engineering*, 38(1):67–86, 2023.
- 802 [He et al., 2022] Kaiming He, Xinlei Chen, Saining Xie,
803 Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
804 toencoders are scalable vision learners. In *Proceedings of
805 the IEEE/CVF conference on computer vision and pattern
806 recognition*, pages 16000–16009, 2022.
- 807 [Hinton and Salakhutdinov, 2006] Geoffrey E Hinton and
808 Ruslan R Salakhutdinov. Reducing the dimensionality of
809 data with neural networks. *science*, 313(5786):504–507,
810 2006.

- 811 [Hsiao *et al.*, 2009] Pei-Yung Hsiao, Chun-Wei Yeh, Shih-
812 Shinh Huang, and Li-Chen Fu. A portable vision-based
813 real-time lane departure warning system: Day and night.
814 *IEEE Transactions on Vehicular Technology*, 58(4):2089–
815 2094, 2009.
- 816 [Li *et al.*, 2019] Xiang Li, Jun Li, Xiaolin Hu, and Jian Yang.
817 Line-cnn: End-to-end traffic line detection with line pro-
818 posal unit. *IEEE Transactions on Intelligent Transporta-*
819 *tion Systems*, 21(1):248–258, 2019.
- 820 [Mehrabi *et al.*, 2022] Ninareh Mehrabi, Fred Morstatter,
821 Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A
822 survey on bias and fairness in machine learning, 2022.
- 823 [Meng *et al.*, 2021] Depu Meng, Xiaokang Chen, Zejia Fan,
824 Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jing-
825 dong Wang. Conditional detr for fast training convergence.
826 In *Proceedings of the IEEE/CVF International Conference*
827 *on Computer Vision (ICCV)*, pages 3651–3660, October
828 2021.
- 829 [Pan *et al.*, 2018] Xingang Pan, Jianping Shi, Ping Luo, Xi-
830 aogang Wang, and Xiaoou Tang. Spatial as deep: Spa-
831 tial cnn for traffic scene understanding. In *Proceedings of*
832 *the AAAI Conference on Artificial Intelligence*, volume 32,
833 2018.
- 834 [Roselli *et al.*, 2019] Drew Roselli, Jeanna Matthews, and
835 Nisha Talagala. Managing bias in ai. In *Companion Pro-*
836 *ceedings of The 2019 World Wide Web Conference, WWW*
837 *'19*, page 539–544, New York, NY, USA, 2019. Associa-
838 tion for Computing Machinery.
- 839 [Tabelini *et al.*, 2021] Lucas Tabelini, Rodrigo Berriel, Thi-
840 ago M Paixao, Claudine Badue, Alberto F De Souza, and
841 Thiago Oliveira-Santos. Keep your eyes on the lane: Real-
842 time attention-guided lane detection. In *Proceedings of*
843 *the IEEE/CVF conference on computer vision and pattern*
844 *recognition*, pages 294–302, 2021.
- 845 [Tang *et al.*, 2021] Jigang Tang, Songbin Li, and Peng Liu.
846 A review of lane detection methods based on deep learn-
847 ing. *Pattern Recognition*, 111:107623, 2021.
- 848 [Vaswani *et al.*, 2023] Ashish Vaswani, Noam Shazeer, Niki
849 Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
850 Lukasz Kaiser, and Illia Polosukhin. Attention is all you
851 need, 2023.
- 852 [Wang *et al.*, 2010] Jyun-Guo Wang, Cheng-Jian Lin, and
853 Shyi-Ming Chen. Applying fuzzy method to vision-based
854 lane detection and departure warning system. *Expert Sys-*
855 *tems with Applications*, 37(1):113–126, 2010.
- 856 [Wang *et al.*, 2011] Xiaoyun Wang, Yongzhong Wang, and
857 Chenglin Wen. Robust lane detection based on gradient-
858 pairs constraint. In *Proceedings of the 30th Chinese Con-*
859 *trol Conference*, pages 3181–3185. IEEE, 2011.
- 860 [Xingang Pan and Tang, 2018] Ping Luo Xiaogang Wang
861 Xingang Pan, Jianping Shi and Xiaoou Tang. Spatial
862 as deep: Spatial cnn for traffic scene understanding. In
863 *AAAI Conference on Artificial Intelligence (AAAI)*, Febru-
864 ary 2018.
- [Yoo *et al.*, 2020] Seungwoo Yoo, Heeseok Lee, Heesoo 865
Myeong, Sungrack Yun, Hyoungwoo Park, Janghoon Cho, 866
and Duck Hoon Kim. End-to-end lane marker detection 867
via row-wise classification, 2020. 868
- [Zhou and Zhou, 2023] Kunyang Zhou and Rui Zhou. End- 869
to-end lane detection with one-to-several transformer, 870
2023. 871
- [Zou *et al.*, 2019] Qin Zou, Hanwen Jiang, Qiyu Dai, Yuan- 872
hao Yue, Long Chen, and Qian Wang. Robust lane de- 873
tection from continuous driving scenes using deep neu- 874
ral networks. *IEEE transactions on vehicular technology*, 875
69(1):41–54, 2019. 876

877 A Methodology section from the midterm

878 The approach we have taken is based on the O2SFormer archi- 935
879 tecture. As mentioned previously, this is one of the most 936
880 recent papers on the topic of lane detection using transform- 937
881 ers and it showcases state-of-the-art performance. Thus, it 938
882 provides a good starting point to experiment with using se- 939
883 quential images in a transformer setting, which is ultimately 940
884 the goal of the project. At this stage, the project is still in 941
885 the implementation phase, therefore, this section of the report 942
886 portrays our progress to this moment and our plans for the 943
887 coming weeks. 944

888 First, we stick to the pre-processing steps that O2SFormer 945
889 does, namely we resize the image into 320x800 image and 946
890 normalize its values. We currently do not do randomized 947
891 horizontal flipping, but we will want to add it in the coming 948
892 weeks. The data is preprocessed so during training it loads 949
893 the inputs and targets in the correct format, and these formats 950
894 are also cached for quicker access. 951

895 The workflow of the system begins with a pre-trained con- 952
896 volutional neural network (ResNet18) - a so-called backbone. 953
897 This CNN is trained on more than a million images from 954
898 the ImageNet database and is commonly used in the litera- 955
899 ture as a feature extractor. ResNet is the CNN used by the 956
900 O2SFormer and we chose ResNet18 specifically as the most 957
901 lightweight ResNet architecture. The data from the CULane 958
902 dataset is fed into the backbone network in batches, the size 959
903 of which we are still experimenting with as we have hardware 960
904 limitations. The output of the backbone per batch is a collec- 961
905 tion of vectors, each with the size of the input image. These 962
906 vectors, which we refer to as segments, contain extracted fea-
907 tures from the images of our input dataset.

908 The next step is to positionally encode the segments com- 963
909 ing from the backbone. Positional encoding is used to pro- 964
910 vide information about the positions of elements in an input 965
911 sequence and is a key component of any transformer-based 966
912 model. Thus, we use positional encoding which computes 967
913 sine and cosine functions in order to preserve the spatial in- 968
914 formation of the sequential images we have. The encoding 969
915 used by us is inspired once again by the O2SFormer. 970

916 Once we have the batches of positionally encoded seg- 971
917 ments, we pass those to the transformer. It naturally consists 972
918 of an encoder and a decoder. At this stage, we make use 973
919 of the PyTorch implementation of these components. The 974
920 O2SFormer paper does not suggest any changes to the en- 975
921 coder part, however, it adds two additional components to the 976
922 decoder - content query and dynamic anchor-based positional 977
923 query. These components are influenced by the Conditional 978
924 DETR ([Meng *et al.*, 2021]) paper which proposes a division 979
925 of the object query in DETR into content query and positional 980
926 query. Object query in DETR is an attention mechanism that 981
927 finds the important features in the encoder output and even- 982
928 tually produces the output of the decoder. This object query 983
929 in DETR and the content and positional queries in the detec- 984
930 tion transformers replace the output of the previous step as an 985
931 input to the decoder. 986

932 To train we have defined our own loss function which com- 987
933 putes the MSE loss between all the points where a lane is 988
934 supposed to be and all the points where it is not supposed to 989
935

be. The final loss then consists of equally-weighted average 935
of both of these sub-losses. We have done this for 2 reasons, 936
first the loss in [Zhou and Zhou, 2023] is not well defined, as 937
the output of the model does not correspond to the target, so 938
some post-processing had to be done, which influences how 939
the gradients for backwards pass are computed. Secondly, we 940
wanted to ensure that the model learns to output something, 941
as there are many more cells with no lane than ones with a 942
lane. 943

Thus far, we have not added a content query and a dy- 944
namic anchor-based positional query as proposed by the 945
O2SFormer. Instead, we feed the output of the previous pre- 946
diction step back into the decoder as normally done in trans- 947
formers for NLP tasks. Additionally, we still experiment with 948
different output formats. Our initial idea was to represent the 949
output as a grid over the image and each cell in this grid to 950
have 3 parameters assigned - angle length and probability. 951
The angle is between the predicted lane and the x-axis of the 952
grid, the length and probability again refer to the lane and 953
how likely it is to have one in that cell of the grid. 954

We have also tried a simpler architecture, where the model 955
directly outputs a mask which is treated as a binary mask in- 956
dicating where the lanes can be found. While in general, pre- 957
dicting the masks directly results in longer training times and 958
the models are easier to overfit, it gives us the opportunity to 959
use the Intersection over Union as our loss function. There- 960
fore it should in theory be actively pushed towards making a 961
better prediction per frame. 962

963 B Evaluation from midterm report

To evaluate the current state of our project, we use the same 964
approach as other researchers working with CULane. The au- 965
thors of the CULane dataset have created separate categories 966
of videos, which correspond to a certain difficulty when try- 967
ing to predict lanes. These categories consist of videos with- 968
out lines, with cars crossing over lines, or videos done at night 969
time. To measure how well does the model predict these lines 970
we use the intersection over union score between binary mask 971
generated from the output of the model and binary mask gen- 972
erated from the ground truth. Be aware that these results were 973
trained for XYZ epochs on only 30% of the training data. 974
These results are currently quite bad, as the models tend to 975
learn to output the entire road rather than the specific lanes. 976
The model that directly predicts a masked image suffers from 977
the same problem. 978

Part of the evaluation is the intensity of computation. Cur- 979
rently we manage to treat 30 frames as one video (which is 980
one batch), which gives an overall window of about 15 sec- 981
onds, as there are 2 frames per second in the dataset. We use 982
30 as there are 180 or 90 frames in each video, and using 983
30 avoids having frames from different videos in one batch. 984
At this stage for both approaches going through 30% of the 985
dataset takes around 1 hour. One epoch over the entire dataset 986
tends to take slightly less than 4 hours. While we can get 987
some preliminary results, we will not be able to compare our 988
models to other models, as they tend to train for around 20 989
epochs, which in our case would take 80 hours. 990

991 **C CULane supplementary material**



Figure 8: Examples for each challenging category in the test set.

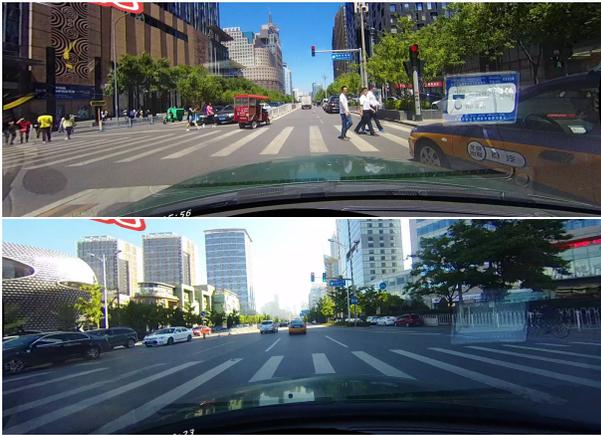


Figure 9: Examples of unannotated lanes in sight.

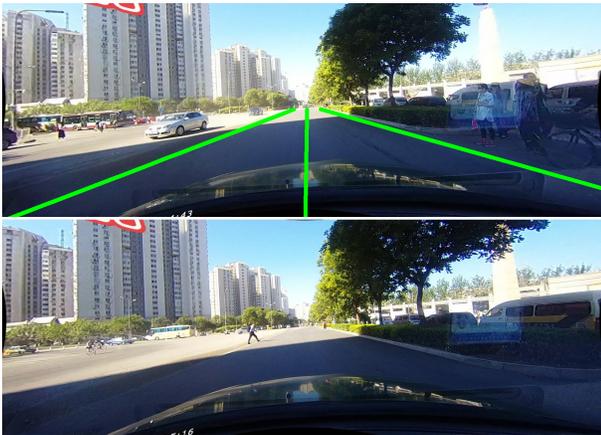


Figure 10: Lane annotations on the first image can be seen to not have much visual indication on the second image.

992 **D Loss formulas and references**

993 **D.1 CE Loss**

994 For more information on the specifics for this loss, and the
 995 Binary CE loss please visit [the pytorch website for CE loss](#)
 996 and [the one for BCE loss](#).

997 **E Model and hardware parameters**

998 For the final version of the model we worked on, we have the
 999 following parameters:

d-model = 2304, nhead=1, out-dim=(224, 224),
 dropout = 0.1, num-decoder-layers = 1,
 batch-size = 32, learning-rate= 1e-5,
 weight-decay = 1e-7.

Hardware used for training:

CPU: AMD Ryzen 9 5900X 12-Core, 11th Gen Intel(R)
 Core(TM) i9-11900KF @ 3.50 GHz, AMD Ryzen 7 5600X

GPU: NVIDIA GeForce RTX 3070 Ti, NVIDIA GeForce
 GTX 1080 Ti

F Reconstruction Results from ResNet34 Auto-Encoder

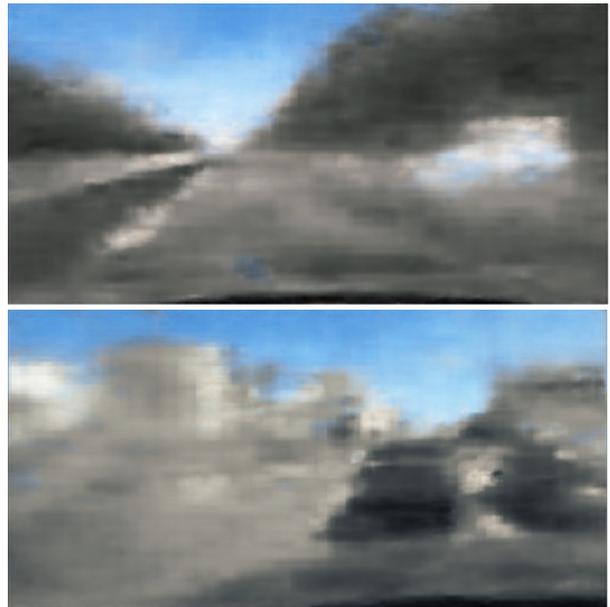


Figure 11: 2 images from the CULane train dataset, reconstructed by a decoder

G Reconstruction from a MAE ViT



Figure 12: Input image to the MAE ViT feature extractor



Figure 13: Output image to the MAE ViT feature extractor



Figure 14: Input image to the MAE ViT feature extractor



Figure 15: Output image to the MAE ViT feature extractor, lanes are clearly visible.

1012 H More results

1013 Figure 17 and Figure 16 present the loss and IoU when training
 1014 without feeding the ground truth to the decoder. We see
 1015 that again both the training and validation losses are decreasing.
 1016 However, we have run this setup for only 6 epochs, starting
 1017 from the checkpoint after the 16th epoch of the previously
 1018 discussed training. Thus, the results are somewhat unstable
 1019 and the validation loss is better than the training one which
 1020 suggests that this model probably finds it easier to generalize,
 1021 but the training set is filled with more instances of outliers.
 1022 Additionally, Figure 19 displays the mean IoU for specific
 1023 frames. One can observe that after approximately 25 epochs
 1024 (first with ground truth being fed as input, then without) the
 1025 network tends to exploit some information about the previous
 1026 frame to get a better prediction. The overall variation is small
 1027 between epochs so we can say that the IoU stays relatively
 1028 stable. When it comes to the test set, Figure 18 illustrates
 1029 the IoU score over different thresholds set for the confidence
 1030 level. We notice that the best results are when the threshold is
 1031 around 0.85 when the mean IoU reaches levels of about 0.12.

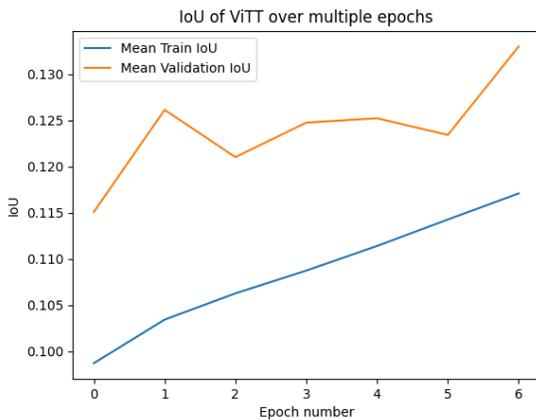


Figure 16: IoU over 6 epochs without using the ground truth in the decoder

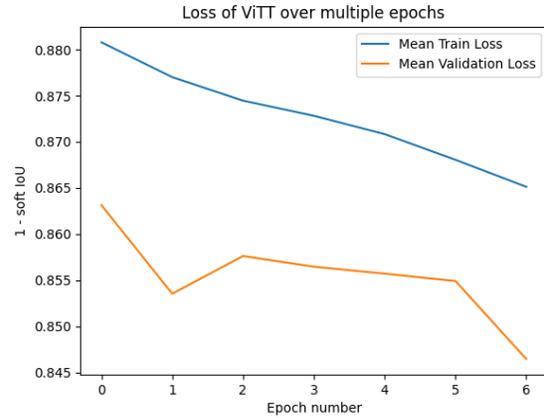


Figure 17: Loss over 6 epochs without using the ground truth in the decode

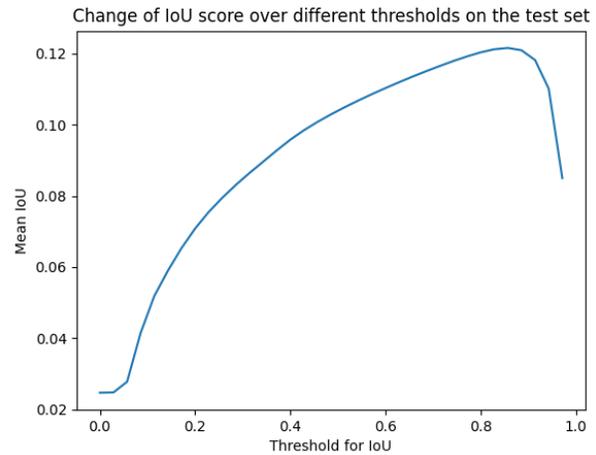


Figure 18: IoU over different thresholds for lane classification.

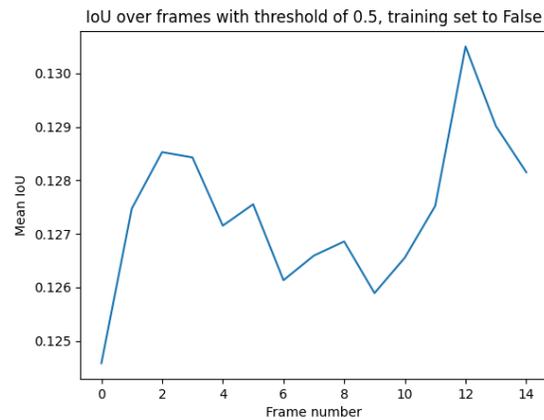


Figure 19: Average IoU for specific frames of the 15 frames which are being consumed by the model.

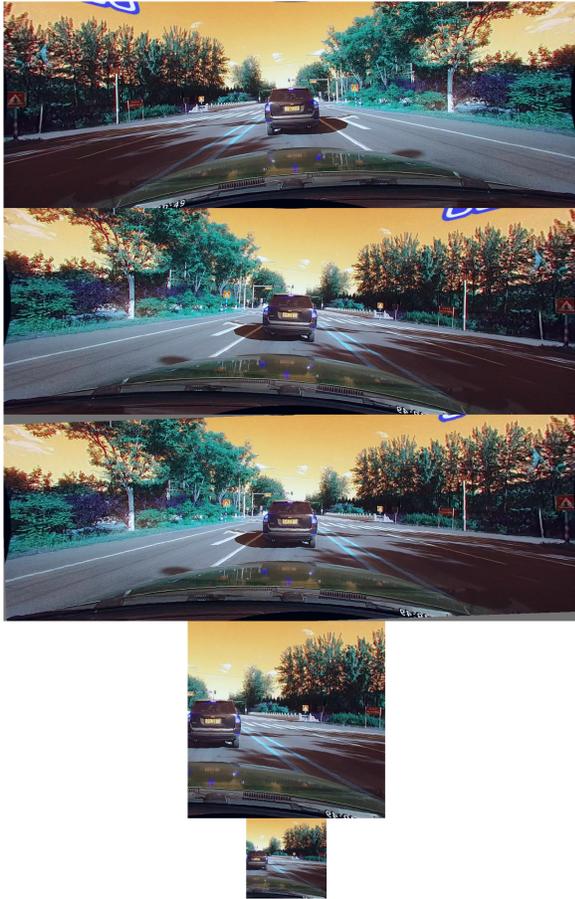


Figure 20: Visualization of the four data augmentation and preprocessing steps in the order they are executed from top to bottom. These are flips, rotations, cropping, and resizing. The colors appear distorted due to pixel value normalization.