

---

# Reproducibility report of *From goals, waypoints & paths to longterm human trajectory forecasting* for ML Reproducibility Challenge 2021

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Scope of Reproducibility

2 The following paper is a reproducibility report for From Goals, Waypoints & Paths To Long Term Human Trajectory  
3 Forecasting [9]. The basic code was made available by the author at *this https url*. We have verified all claims and  
4 results from the experiments mentioned in the paper to support the claims. The central claim of YNet is that it sets  
5 state-of-the-art short and long-term prediction standards by a multi-modal network employing both segmentation  
6 matrices and past trajectory heat-maps together.

## 7 Methodology

8 The model essentially combines a segmentation map and past trajectory heatmaps to encode a combined input to three  
9 sub-networks modelled after the U-Net architecture [12]. The author's code was used to benchmark the claims, and  
10 some experiments were performed thereafter. Free-to-use platforms like Google Colaboratory and Kaggle were used to  
11 train these models. We have reproduced the code base in PyTorch Lightning (originally in PyTorch Ignite) and found  
12 consistent results across the board.

## 13 Results

14 Through our testing, we were able to come within 2% of the proposed metrics on certain datasets like Stanford Drone  
15 Dataset (SDD) and ETH/UCY, implying the author's claims are sanguine and reproducible on varied hardware. However,  
16 certain results such as long-term predictions on the Intersection Drone (InD) dataset were quite different; the probable  
17 reasons for which have been discussed.

## 18 What was easy

19 Obtaining the proposed results on the SDD and InD datasets was easy. Well-documented interactive notebooks for  
20 training and testing along with the requisite data for SDD were provided with the codebase. The code could be run with  
21 minimal changes overall.

## 22 What was difficult

23 The codebase and data provided by the authors were incomplete, and contained various redundancies and unused  
24 methods, making it difficult to follow. The requisite code and data required to reproduce the experiments on the  
25 ETH and UCY datasets were completely missing. These factors were compounded by stringent computational power  
26 requirements, which were difficult to fulfill for students without access to server-grade computation.

## 27 Communication with original authors

28 Attempts were made to contact the authors regarding some doubts, which went without response. Running the  
29 experiments thereafter were done based on our understanding of the paper and code.

## 30 **1 Introduction**

31 The paper reproduced in this report aims to tackle multiple pedestrian trajectory predictions using rich multi-modal  
32 predictions for the use of autonomous vehicles, social robots, etc. Earlier approaches to this problem have been  
33 auto-regressive in nature [1][8][14], i.e., using  $n$  points (or, analogically, data from the last  $t$  seconds) from the dataset  
34 to produce the next immediate point, and recurring this process.

35 In this paper, the trajectory distribution, viz. the path taken by a pedestrian is conceived to have been influenced majorly  
36 by two factors:

- 37 • Epistemic: The conscious will of the pedestrian to reach a particular goal.
- 38 • Aleatoric: The unknown and unexpected changes in the environment influencing the path they take to reach  
39 the goal.

40 The proposed architecture incorporates this multi-modality. An explicit probability distribution of the many possible  
41 broad future trajectories is predicted first (modelling the *where* of the agent). Then, random future points of the trajectory  
42 are taken in conjunction with the sampled way-points to obtain probability maps over the remaining predicted points  
43 (modelling the *how* of the agent).

44 To formulate this report, we have experimented on the author's code by adding/removing social pooling layers and  
45 employing visualisation tools. We have tried the unique idea of multi-dataset training wherein we train the model for  
46 long on a particular dataset, and then immediately introduce it to a completely new dataset. We also performed some  
47 experiments such as shifting the prediction origin to different previously predicted points instead of the one closest in  
48 time to the present. These experiments are explained in detail in the following sections.

## 49 **2 Scope of reproducibility**

50 The problem of multi-modal trajectory prediction is key to unlocking vehicular intelligence and autonomous navigation.  
51 The problem this particular paper aims to address is finding pedestrian trajectories in an environment crowded with  
52 similar and/ or different interacting agents. By extension, the scope of using such architecture is beyond pedestrians, as  
53 virtually any human or non-human agent navigating crowded terrains that may benefit from segmentation may employ  
54 such mechanics of trajectory prediction.

55 The central claims of the paper can be summarized as follows:

- 56 • Conditioned way-point predictions: The model performs better than previous works as trajectories are  
57 conditioned in a two-stage manner, with aleatoric predictions conditioned upon epistemic ones. This provides  
58 stricter constraints on the final set of predictions by modelling them explicitly, as opposed to SGAN [5],  
59 SoPhie [13] and other attention based mechanisms that produce a diverse set of trajectories.
- 60 • Scene Segmentation: The model performs better than contemporary models as semantic information about the  
61 scene is accounted for. The paper considers as input, both the segmentation map and trajectory heatmap of  
62 probabilities. The segmentation step is a novel addition that classifies the possible trajectory avenues of the  
63 pedestrian. Intuitively, this can be thought of as follows: Given a predicted valid goal of the pedestrian, he is  
64 highly unlikely to climb a wall to achieve it. Rather, he shall traverse his current course (say, a park track).  
65 The segmentation steps performs better than previous non-segmented attention mechanisms.
- 66 • Long Term Prediction: The model uses these techniques to achieve significantly improved results on long  
67 horizon trajectory prediction as well as short horizon.

## 68 **3 Methodology**

69 We used the GitHub repository provided by the author as the base. However, it only contained the base model for results  
70 on the different data sets. In order to reproduce the rest of the experiments, we had to make changes accordingly.

71 **3.1 Model descriptions**

72 The problem of multi-modal trajectory prediction can be formulated as prediction of future trajectory given past  
73 positions of pedestrians in the scene. This section has been referenced from the original paper [9] (Section 3).

74 The scene image is first processed by a segmentation network, producing segmentation map  $S$  (dividing the image in  
75 various classes) of the same spatial size as image. In a parallel branch the past trajectories are embedded in a trajectory  
76 heatmap. Concatenation of both produces the heatmap tensor  $H_s$ . For each frame  $n$  in the input, the heatmap is  
77 calculated as

$$\mathbf{H}(n, i, j) = 2 \frac{\|(i, j) - \mathbf{u}_n\|}{\max_{(x,y) \in \mathcal{I}} \|(x, y) - \mathbf{u}_n\|} \quad (1)$$

78 The heatmap and semantic maps are concatenated and fed into the encoder branch  $U_e$  of the network.

79 The subsequent architecture consists of 3 sub-networks  $U_e$ ,  $U_g$  and  $U_t$ .  $U_e$  which is an encoder like U-Net [12] is used  
80 in the model architecture to process the tensor  $H_s$ . It has a total of  $N_{U_e}$  blocks, it reduces the dimensions of the  $H \times W$   
81 to  $H_U \times W_U$  and increases the channel depth. The final representation is termed as  $H_{U_e}$  which is then passed in the  
82 goal decoder  $U_g$  and the trajectory decoder  $U_t$ .

83 The next step is termed as the "Goal & Waypoint Heatmap Decoder" in which the output maps of  $U_e$  at various spatial  
84 resolutions are passed in  $U_g$  which is modelled from U-Net. The output is passed through a deconvolution layer, which  
85 involves the application of a transpose convolution, effectively expanding the previous feature map, spatially doubling  
86 the resolution in every block. The encoder map from the respectively sized input layer is concatenated. To attain  
87 the final resolution of the goal, heatmap feature merging is done. Therefore, it can be said a U-Net block consists of  
88 deconvolution, feature merging and convolution layers.

The final branch  $U_t$  is termed as the "Trajectory Heatmap Decoder". The waypoint distributions from  $U_g$  are sampled  
using the softargmax operation

$$\text{softargmax}(X) = \left( \sum_i i \frac{\sum_j e^{X_{ij}}}{\sum_{i,j} e^{X_{ij}}}, \sum_j j \frac{\sum_i e^{X_{ij}}}{\sum_{i,j} e^{X_{ij}}} \right)$$

89 A heatmap tensor  $\mathbf{H}_{U_g}$  is generated using these samples. Each heatmap is downsampled to its corresponding size from  
90 the architecture. These heatmaps are concatenated with the respective blocks from  $\mathbf{H}_{U_e}$  which goes through  $U_t$  for a  
91 decoding phase.

92 **3.2 Datasets**

93 All annotations were preprocessed to match the format ['trackId', 'frame', 'x', 'y', 'sceneId',  
94 'metaId']. For experiments based on varying the prediction window, the data was preprocessed with different  
95 trajectory lengths.

96  
97 **Stanford Drone Dataset (SDD):** [11] The dataset by default contains annotations for 10,300 unique agents across 6  
98 classes, of which 5232 belong to the class of pedestrians. Trajectories are sampled at FPS = 30 in 2D image coordinates.  
99 We use the pre-processed data provided by the authors, which has been downsampled to FPS = 2.5 for short term  
100 training and FPS = 1 for long term. The lengths of the input sequences  $n_p$  are 8 and 5 respectively, while the those  
101 of the output  $n_f$  are 12 and 60 respectively. All trajectories not belonging to the pedestrian class or of insufficient  
102 length ( $< n_f + n_p$ ) are dropped. The midpoints of the bounding boxes are considered to be the ground truth positions.  
103 Trajectories are split at temporal discontinuities and a staggered sliding window is used to split long trajectories. The  
104 resultant is a set of 1502 trajectories. A semantic map with 5 classes is generated. There is a train/test split of 30 scenes  
105 for training and 17 for testing.

106  
107 **Intersection Drone Dataset (InD):** [3] The dataset by default contains 11,500 trajectories across 3 classes, in 32 scenes  
108 at 4 distinct locations. Trajectories are sampled at FPS = 25 in 2D world coordinates. We perform the preprocessing  
109 described in the paper, which involves downsampling the data to FPS = 1 for  $n_p = 5$  and  $n_f = 30$ , filtering out  
110 non-pedestrians, filtering out short ( $< n_f + n_p$ ) trajectories, splitting long (using the sliding window technique) and  
111 discontinuous trajectories. The coordinates are brought into image coordinates by using the scale factors and cropping

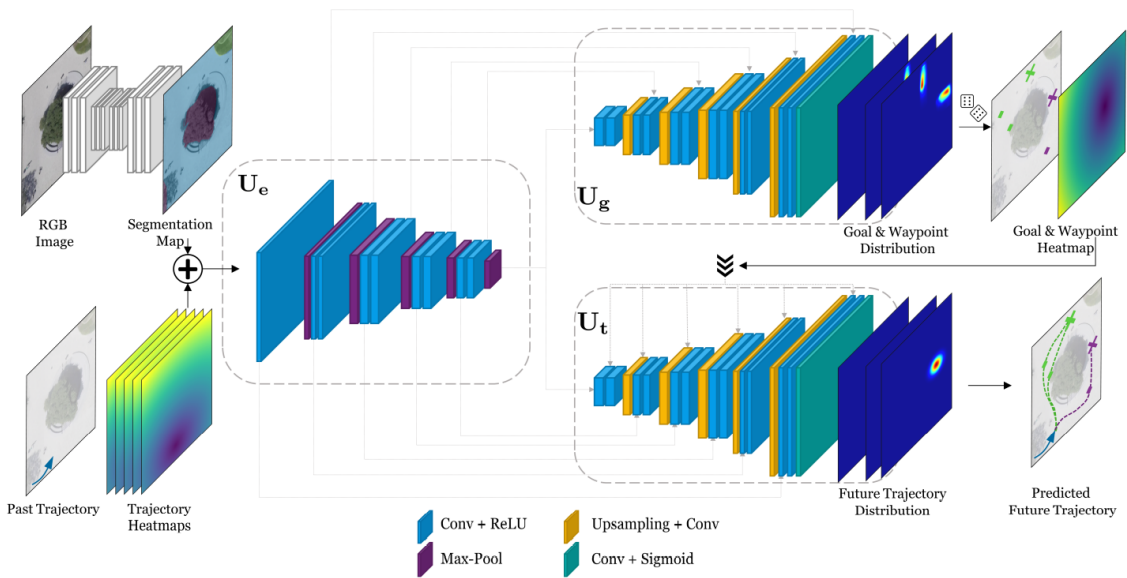


Figure 1: Graphical view of model architecture (referenced from parent paper)

112 parameters provided in the paper as cited. The resultant is a set of 1396 trajectories. The scenes of one location (ID 4)  
 113 are used for testing while those of the remaining 3 are used for training.

114

115 **ETH and UCY datasets (ETH/UCY):** Combined, the dataset contains trajectories for 1536 pedestrians, in 9 scenes at 5  
 116 distinct locations. Trajectories are sampled at FPS = 2.5 in 2D world coordinates. The authors claim to use preprocessed  
 117 data from [5], however this is not usable on account of being normalized with unknown parameters and being in the  
 118 incorrect format. We instead use preprocessed data provided by the authors in response to an issue raised on the GitHub  
 119 repository. We take  $n_p = 8$  and  $n_f = 12$  as per the paper. We use the homography matrices provided with the datasets  
 120 to transform the pixels into world coordinates. A leave-one-out cross-validation strategy is employed.

### 121 3.3 Hyperparameters

122 Several hyperparameters were experimented with in this paper. Those of particular importance are the following:

- 123 •  $K_e$ : The model aims to produce  $K_e$  possible future trajectories due to the epistemic uncertainty of final goal.  
 124 This is a hyperparameter that can be tuned to find the optimal value for any given probability map as input.
- 125 •  $K_a$ : After the end-point prediction distribution, the path(s) taken to reach the most likely of them constitutes  
 126 absolute randomness when not conditioned on aleatory factors and environmental interactions. Thus, given the  
 127 goal, the model produces  $K_a$  separate predictions for path.
- 128 •  $T$ : The temperature parameter T during sampling can be visualised as a scaling factor for the generated  
 129 heat-map. It is used to control the trade-off between diversity and precision; a lower value meaning the  
 130 predictions are condensed in a smaller spatial density and vice versa. Intuitively, a higher value of T should be  
 131 used for long-term predictions, but this parameter can still be tuned to gauge the power of the model.

132 All hyper parameters were tuned by random searches and heuristic guesses instead of brute/ grid searches or Bayesian  
 133 techniques, mainly due to constraints posed by very high computational resource requirements. However, the trends in  
 134 accuracy could still be predicted and reasoned as the effects of changing the values were both experimentally visible  
 135 and logically explainable. These points have further been discussed in the Results and Discussions sections.

Hyperparameters	Value
$K_e$	5 or 20
$K_a$	1
$T$	1.8
Number of epochs	100
Batch Size	8 (4 on long term)
Learning Rate	1E-4 (optimal)
Semantic Classes	6 (3 for ETH/UCY)
Waypoints	11

Table 1: Hyperparameters used in the paper

### 136 3.4 Experimental setup and code

137 The code for this experiment is set-up mainly in the `train.py`, `evaluate.py` and `test.py` Python files that import  
 138 helper methods defined in python files in the `utils` folder. Python notebooks are given to facilitate running different parts  
 139 of the code. Detailed instructions about each, the presence of pre-trained weights and/ or pre-processed files and other  
 140 relevant information is given in the ReadMe section of the repository.

141 The main metrics of interest are the ADE (mean L2-norm distance between all future ground truth and predicted  
 142 points) and FDE (mean L2-norm distance between final future ground truth and predicted points). The accuracy of all  
 143 experiments are validated with these metrics, where a lower value means a more accurate result.

### 144 3.5 Computational requirements

145 All experiments were run using Google Colaboratory, whose back-end has the Tesla P100 GPU. The technical  
 146 specification of the GPU is that it has 3584 CUDA cores, 16GB CRAM and a 4096-bit memory interface.

147 The run-times we faced on such a setup for the different experiments is quite long. For example, running the code  
 148 without any ablations on the SDD dataset took roughly 10 minutes for a single epoch. Of course, this value is dependant  
 149 on other hyperparameters such as batch size.

## 150 4 Results

151 The results we obtained are listed below. Upon comparison, we are confident that the results resemble those in the paper.  
 152 However, due to the lack of extensive computational capabilities, we were forced to limit our training to a fraction  
 153 of what was done in the original paper. Despite this, we have deeply analysed fitting and convergence trends and are  
 154 confident that the model does at least as well as claimed, and even better in some experiments.

155 All results were logged with ease with the WandB solution [2]. In general, extensive overview of error trends could be  
 156 gauged from auto-generated graphs, which cemented our beliefs of convergence and correctness.

### 157 4.1 Results reproducing original paper

#### 158 4.1.1 Performance of model as compared to baselines

159 Our reproduced model functions better than all previous baselines, and satisfactorily close to the results of YNet as cited  
 160 in the paper.

	K=20					K=5			
	P2TIRL[4]	SimAug[7]	PECNet[10]	Y-net (Paper's)	Ours'	TNT[15]	PECNet[10]	Y-net (Paper's)	Ours'
ADE	12.58	10.27	9.96	7.85	<b>8.85</b>	12.23	12.79	11.49	<b>12.36</b>
FDE	22.07	19.71	15.88	11.85	<b>12.23</b>	21.16	29.58	20.23	<b>20.18</b>

Table 2: Short temporal horizon forecasting results on SDD

161 **4.1.2 Performance of model for different datasets of ETH/UCY: Importance of social masking**

162 This table is produced separately because it addresses the importance of social masking. The paper results are with  
 163 masking, while ours are without.

	ADE		FDE	
	YNet (Paper's)	Our's	YNet (Paper's)	Our's
ETH	0.28	<b>0.38</b>	0.33	<b>0.60</b>
HOTEL	0.10	<b>0.69</b>	0.14	<b>0.97</b>
UNIV	0.24	<b>0.35</b>	0.41	<b>0.61</b>
ZARA1	0.17	<b>0.31</b>	0.27	<b>0.423</b>
ZARA2	0.13	<b>0.34</b>	0.22	<b>0.57</b>

Table 3: Short temporal horizon forecasting results on several datasets of ETH/UCY without social masking.

164 **4.1.3 Turning TTST and CWS on and off**

165 TTST and CWS are heuristics designed to improve sampling. TTST encourages clustering samples in high-density  
 166 regions by roughly thresholding and clustering the probability distribution. CWS discourages sampling erratic  
 167 trajectories by assuming points to be sampled between two known points roughly divide the line segment joining them.  
 168 The roughness is modelled using Gaussian distributions. We experiment with various possible state to verify their effect  
 169 on error.

	SDD (paper's / ours)			IND (paper's / ours)	
	x	x	✓	x	✓
TTST	x	x	✓	x	✓
CWS	x	✓	✓	✓	✓
ADE	65.00 / <b>46.52</b>	52.31 / <b>46.67</b>	47.94 / <b>46.52</b>	17.77 / <b>4.85</b>	14.99 / <b>4.90</b>
FDE	86.98 / <b>62.23</b>	86.98 / <b>63.23</b>	66.71 / <b>62.23</b>	28.52 / <b>8.85</b>	21.13 / <b>9.23</b>

Table 4: Effect of TTST and CWS on SDD on inD

170 **4.1.4 Hyperparameter tuning -  $K_a$ ,  $K_e$  and  $T$**

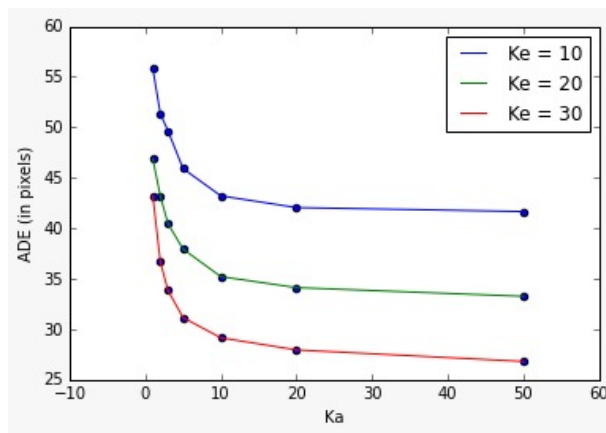


Figure 2: Variation of ADE with grid search wise changes in  $K_a$ ,  $K_e$

171 **4.2 Results beyond original paper**

172 **4.2.1 Generalization**

173 One of our main findings beyond the paper was the generalizing power of the model, abstracted by its potential to be  
 174 used as a transfer learning model. Given that the data can be processed in a similar manner outside the domain of the

175 actual model, we observed much-improved results when trained for a very short time on a completely new dataset. To  
 176 explore this further the idea of Fine-tuning was explored in which once the Y-net model was trained on Dataset A, the  
 177 final weights were considered as the pretrained weights for a new training and the model was further trained on Dataset  
 178 B for very few epochs.

179 In this way the model not only remembered the previous training features but also adapted the conditions for the new  
 180 dataset. This method proved to improve the performance of the model and is computationally very inexpensive.

Model Trained on	Epochs	Model tested on	No. of Epochs further trained for	ADE	FDE
inD	300	inD	0	14.99	21.13
SDD longterm	300	inD	0	10.67	17.21
SDD longterm	300	inD	1	5.032	8.767
SDD longterm	300	inD	2	4.967	8.844
SDD longterm	300	inD	3	4.822	8.699
SDD longterm	300	inD	<b>4</b>	<b>4.59</b>	<b>8.090</b>

Table 5: Long term trajectory forecasting Results on transferred dataset

## 181 5 Visualisations

182 Some real world trajectories on actual reference images are provided below. The lines are enlarged for clarity. It can  
 183 clearly be seen that the model works fantastically well in real-life scenarios to predict trajectories. Segmentation has  
 184 worked well in these cases, with no class overlap except cases when the trajectory itself goes across two different  
 185 semantic classes like in figure (b).

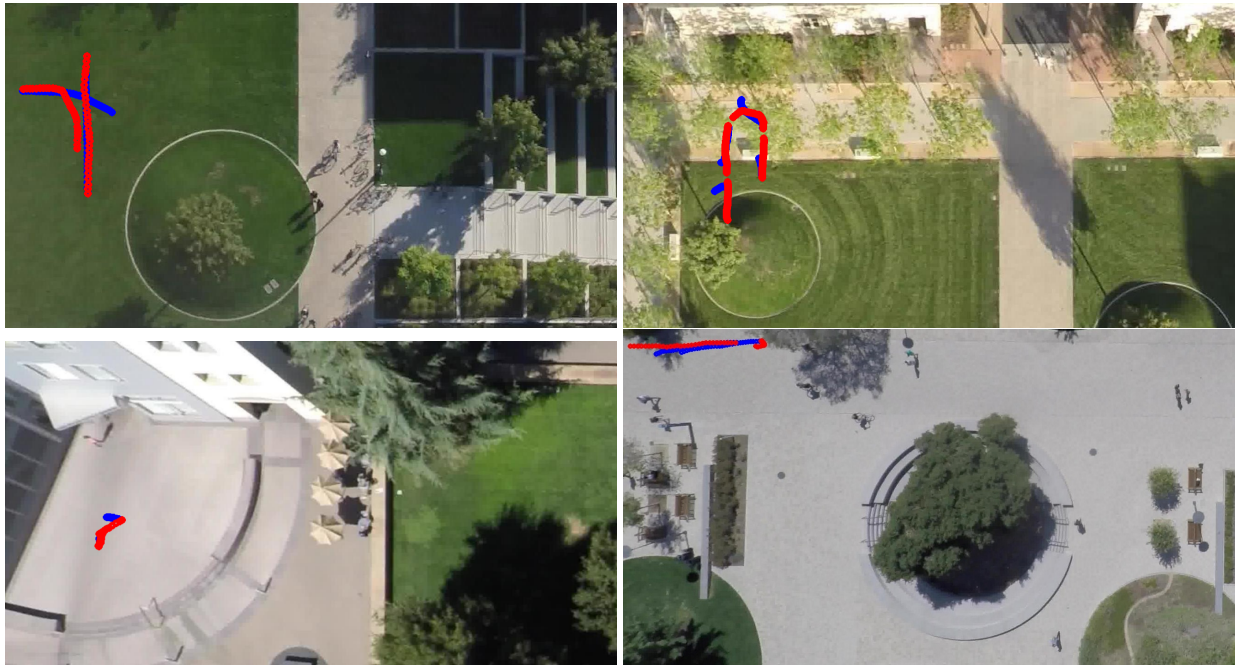


Figure 3: SDD ground truth (red) and predictions (blue) on (a) hyang-0-1 (b) hyang-1-1 (c) gates-2-2 (d) coupa-1-3 reference images. Figure (c) is a short-term prediction with minimal movement, while the other three are long-term with temporal horizon of 30 seconds.

## 186 6 Discussion

187 Many obstacles were faced in the reproduction of the results, particularly for students with limited access to server grade  
 188 computational resources. The codebase and data had many redundancies and omissions, requiring some experiments to

189 be recreated from scratch. Despite these challenges, our experiments achieve a satisfactory reproduction of the paper.  
190 However, some discrepancies were observed. Our results on the InD dataset were significantly better than those cited in  
191 the paper. This may be due to model optimizations in PyTorch Lightning or fortunate random initialization of weights.  
192 The experiments on the SDD and ETH/UCY datasets were found to be consistent with the paper. The predicted  
193 trajectories represented state-of-the-art accuracy, even more so with the TTST and CWS sampling techniques enabled.  
194 We observed enhanced accuracy with dataset dilution and marginal training on a new dataset. The long-term prediction  
195 results were viable, viz. significantly better than contemporary models, enabling this model to be used in a real-time  
196 prediction stack for trajectory prediction.

## 197 **6.1 Further discussion on the results**

198 Table 4 confirms that the usage of TTST and CWS markedly improves the accuracy of the final results by increasing  
199 the tendency to draw samples from relevant points in the probability distribution. However, this comes at the cost of  
200 increased computational complexity.

201 The model is robust towards changes in context and hence can be extended to a wide variety of applications, as  
202 evidenced by Table 5. Good transfer performance indicates the architecture is the dominant factor in our results as  
203 opposed to extraneous factors such as sampling techniques, demonstrating its strength. This experiment also highlights  
204 the potential of transfer learning in improving neural network performance, especially for complex models like this,  
205 which reap the benefit of carrying over a dense field of features.

206 The crux of the paper, the predictive power of the chosen multi-modality, is succinctly demonstrated by Figure 2.  
207 We see a marked improvement in inference as we increase both  $K_a$  and  $K_e$  independently of each other. This direct  
208 relationship indicates that the approach of sequentially predicting epistemic and aleatoric distributions is significant,  
209 and verifies this paper’s contribution to the state-of-the-art of pedestrian trajectory prediction.

210 There are significant increases in all errors upon removing social masking and pooling as seen in table 3. This is a  
211 central claim of a paper, i.e. aleatory interactions from the surroundings are a crucial factor in determining the best path  
212 taken. Modelling them using the segmentation ResNet-101 [6] is evidently better than using deterministic criteria to  
213 model these interactions.

## 214 **6.2 What was easy**

215 The experiments on the SDD and InD datasets required minimal effort to reproduce. The authors provide interactive  
216 Python notebooks for training and testing, along with all the requisite scripts and most of the data to run them. Due to  
217 the modularity of the code, performing the ablation study was also easy.

## 218 **6.3 What was difficult**

219 There were significant challenges faced in this reproduction. Due to limited computational resources, training the  
220 extensive CNN was problematic (mainly owing to a per-epoch training time of 30-60 minutes, despite a small batch  
221 size of 4). Further, the code, data, pre-trained weights, semantic maps, semantic models for the experiments on  
222 the ETH/UCY datasets were missing from the repository, rendering it impossible to exactly reproduce the authors’  
223 experiments. The paper suggests preprocessed data from [5] be used, however we found it was unusable as that data had  
224 been normalized with unknown parameters. The preprocessing functions provided by the authors for ETH/UCY could  
225 not be used as it was not possible to fulfill some arguments. There was an error in the pre-trained weights provided for  
226 the long term SDD experiment, which caused a tensor dimension mismatch during testing.

227 The codebase contained some unused methods. There were some redundant parameters, for example `batch_size = 4`  
228 in the yaml configuration file and `BATCH_SIZE = 8` declared in the training notebooks. Some lines of code were  
229 commented-out without documentation. These factors made it difficult to follow and debug the code.

## 230 **6.4 Communication with original authors**

231 Multiple attempts to contact the authors were made over a 2 month period. Some doubts with the paper and the absence  
232 of ETH/UCY experiments from the codebase were raised. However, there was no response from the authors.



233 **References**

- 234 [1] Apratim Bhattacharyya, Michael Hanselmann, Mario Fritz, Bernt Schiele, and Christoph-Nikolas Straehle.  
235 Conditional flow variational autoencoders for structured sequence prediction, 2020.
- 236 [2] Lukas Biewald. Experiment tracking with weights and biases, 2020. Software available from wandb.com.
- 237 [3] Julian Bock, Robert Krajewski, Tobias Moers, Steffen Runde, Lennart Vater, and Lutz Eckstein. The ind dataset:  
238 A drone dataset of naturalistic road user trajectories at german intersections. In *2020 IEEE Intelligent Vehicles*  
239 *Symposium (IV)*, pages 1929–1934, 2020.
- 240 [4] Nachiket Deo and Mohan M. Trivedi. Trajectory forecasts in unknown environments conditioned on grid-based  
241 plans, 2021.
- 242 [5] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social GAN: socially acceptable  
243 trajectories with generative adversarial networks. *CoRR*, abs/1803.10892, 2018.
- 244 [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*,  
245 abs/1512.03385, 2015.
- 246 [7] Junwei Liang, Lu Jiang, and Alexander Hauptmann. Simaug: Learning robust representations from simulation for  
247 trajectory prediction, 2020.
- 248 [8] Karttikeya Mangalam, Ehsan Adeli, Kuan-Hui Lee, Adrien Gaidon, and Juan Carlos Niebles. Disentangling  
249 human dynamics for pedestrian locomotion forecasting with noisy supervision, 2020.
- 250 [9] Karttikeya Mangalam, Yang An, Harshayu Girase, and Jitendra Malik. From goals, waypoints & paths to long  
251 term human trajectory forecasting. In *Proc. International Conference on Computer Vision (ICCV)*, October 2021.
- 252 [10] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and  
253 Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned trajectory prediction, 2020.
- 254 [11] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Learning social etiquette: Human  
255 trajectory understanding in crowded scenes. In *ECCV*, 2016.
- 256 [12] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image  
257 segmentation. *CoRR*, abs/1505.04597, 2015.
- 258 [13] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, and Silvio Savarese. Sophie: An attentive GAN  
259 for predicting paths compliant to social and physical constraints. *CoRR*, abs/1806.01482, 2018.
- 260 [14] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible  
261 trajectory forecasting with heterogeneous data, 2021.
- 262 [15] Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Benjamin Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen,  
263 Yuning Chai, Cordelia Schmid, Congcong Li, and Dragomir Anguelov. Tnt: Target-driven trajectory prediction,  
264 2020.