

31 after the sixth way-point (and later) were larger than the claimed results. Similar fluctuations were
32 observed for FDE as well, but the relative trends support the paper’s claim.

33 **Communication with original authors**

34 We have not contacted any of the original authors as all the results were reproduced satisfactorily.

35 **1 Introduction**

36 The paper reproduced in this report aims to tackle multiple pedestrian trajectory predictions using rich
37 multi-modal predictions for the use of autonomous vehicles, social robots, etc. Earlier approaches
38 to this problem have been auto-regressive in nature, i.e., using n points (or analogically, data from
39 the last t seconds) from the dataset to produce the immediately next point, and then this process is
40 recurred.

41 In this paper, the end-point distribution conditioned on the past trajectory and the past trajectory
42 features are modeled separately for each pedestrian. The future trajectory points are predicted based
43 on the past and features from other pedestrians via social pooling. An assumption in this model is the
44 absence of passive pedestrians or the fact that each pedestrian has an actual preconceived end-point
45 or destination and is motivated to reach it.

46 To formulate this report, we have experimented on the author’s code by adding/removing social
47 pooling layers, using truncation tricks, visualisation tools, and changing between CVAE and VAE
48 architectures to verify all the claims made by the author described in detail below. We also performed
49 some experiments such as shifting origin to the current point, using different architecture for encoder
50 and decoder networks with the hope of improving the results, which are also described in detail at the
51 end.

52 **2 Scope of reproducibility**

53 The paper revolves around the claim that an important component of predicting the trajectory is
54 the destination in multi trajectory forecasting. If the destination for the pedestrian is clear, then
55 the trajectory can be easily resolved using a separate network that takes the past trajectory and the
56 destination as input taking into account social interactions among fellow pedestrians. Hence the
57 central idea and claim of the paper is to use Conditional Variational Auto Encoder (CVAE) to get the
58 latent variable encoding conditioned on the destination from the ground truth, use the latent variable
59 to infer the predicted destination, and use it for predicting the rest of the future trajectory. We take
60 k samples of the latent variable for testing purposes to predict k different admissible trajectories as
61 output for different destinations derived from the latent encoding. The overall reduction in the value
62 of best ADE and FDE values for the Stanford Drone, ETH/UCY datasets by using the CVAE network
63 is the central claim of the paper.

64 To support the argument that indeed given the destination, the rest of the predicted trajectory con-
65 tributes much less error than the previous state of the art methods such as SGAN, which directly
66 predict the future trajectory, the paper performs an ablation study where they give the ground truth
67 of a way-point which they call as oracle instead of the best one from taking k samples of the latent
68 variable to get the decoupled error of predicting the trajectory. The results strongly support the
69 argument.

70 Further, they also experimented with different values of k to show that FDE tends to 0 as k increases
71 and ADE tends to a certain value, which also shows the decoupled error in predicting the rest of the
72 trajectory.

73 This paper also introduces a non-local social pooling layer and a “truncation-trick,” which improves
74 diversity and multi-modal trajectory prediction performance.

75 Hence the claims can be summarized as follows:-

- 76 1. Conditioning the destination on the past trajectory using CVAE helps in explicit decoupling
77 of the destination prediction and path prediction errors. It hence helps reduce the destination
78 prediction error and the subsequent path prediction error.

- 79 2. Using the social pooling layer helps reduce the error in predicting the path given the history
80 and the destination.
- 81 3. Using truncation trick i.e., truncating the distribution for fewer values of k from which
82 samples are taken helps reduce the destination prediction error. Also, taking a higher sigma
83 value for larger values of k reduces the error.

84 3 Methodology

85 We used the GitHub repository provided by the author as the base. However, it only contained the
86 base model for results on the drone dataset. In order to reproduce the rest of the experiments, we had
87 to make changes accordingly.

88 3.1 Model descriptions

89 The model used in the paper consists of 2 parts:

90 First, the CVAE or Conditional Variational AutoEncoder to get the representation of the latent variable
91 conditioned on destination and given the past trajectory.

92 Second, the predictor network consists of social pooling layers and an MLP network to get the future
93 trajectory.

94 A representative diagram of the network is given in figure 1 and the architecture parameters are shown
95 in table 1.

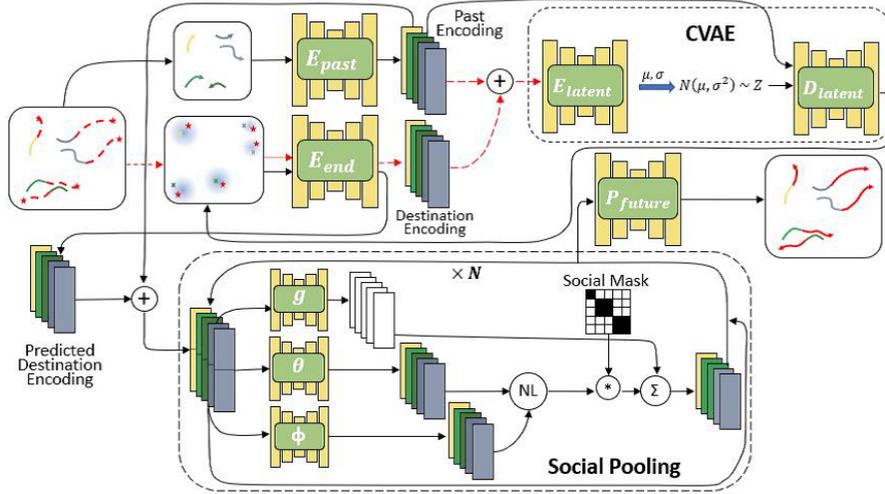


Figure 1: Model architecture

$$ADE = \frac{\sum_{j=t_i+1}^{t_p+t_f+1} \|\hat{\mathbf{u}}_j - \mathbf{u}_j\|_2}{t_f} \quad (1)$$

$$FDE = \|\hat{\mathbf{u}}_{t_p+t_f+1} - \mathbf{u}_{t_p+t_f+1}\|_2 \quad (2)$$

$$\mathcal{L} = \lambda_1 \underbrace{D_{KL}(\mathcal{N}(\mu, \sigma) \parallel \mathcal{X}(0, \mathbf{I}))}_{KL \text{ Div in latent space}} + \lambda_2 \left(\underbrace{\|\hat{\mathcal{G}}_c - \mathcal{G}_c\|_2}_{AEL} + \underbrace{\|\hat{\mathcal{T}}_f - \mathcal{T}_f\|_2}_{ATL} \right)^2 \quad (3)$$

96 **3.2 Datasets**

97 We used Stanford Drone [5] and ETH [4] / UCY [2] datasets. The Stanford drone dataset was given
 98 in the author’s code, but ETH/UCY was not given, so we took the dataset from opensource source.

	Network Architecture
E_{way}	2 -> 8 -> 16 -> 16
E_{past}	16 -> 512 -> 256 -> 16
E_{latent}	32 -> 8 -> 50 -> 32
D_{latent}	32 -> 1024 -> 512 -> 1024 -> 2
θ, Φ	32 -> 512 -> 64 -> 128
g	32 -> 512 -> 64 -> 32
$P_{predict}$	32 -> 1024 -> 512 -> 256 -> 22

Table 1: Model Architecture

99 **3.3 Hyperparameters**

100 We used Hyperparameters given in the paper. We occasionally changed them accordingly to perform
 101 the ablation studies described below.

102 **3.4 Experimental setup**

103 We ran code in google colab with GPU (NVIDIA-SMI 450.36.06 Driver Version: 418.67 CUDA
 104 Version: 10.1).

105 **3.5 Computational requirements**

106 Typically, it took less than an hour to train the model both for the drone and ETH/UCY datasets.

107 **4 Results**

108 The following experiments/ablation studies support the claims made earlier. A detailed description of
 109 the experiments and their results to support the claim are listed below:-

110 **4.1 Experiment on drone dataset (with and without social pooling, truncation trick)**

111 Stanford drone dataset: We did it with social pooling and got results within 95% accuracy from claim
 112 results. The preprocessed dataset for train and test were given on GitHub (by author). We used them
 113 to verify the results. We did two experiments with n-samples 5 and another with n-samples 20 as
 114 required for reproducing the results in the first table of the paper.

	O-S-TT	O-TT	Ours	PECNet-Ours
K	20	20	5	20
ADE	10.56 / 10.47	10.23 / 10.19	12.79 / 14.16	9.96/10.04
FDE	16.72 / 16.43	16.29 / 15.9	25.88 / 26.73	15.96/16.20

Table 2: Comparisons of our results against those of the authors’ and previous state-of-the-art methods. -S’ ‘-TT’ represents ablations of our method without social pooling truncation trick. We report results for in pixels for both K = 5 20 and for several other values of K. The format for each cell is <claimed result> / <reproduced result>

115 **4.2 Experiment on ETH/UCY datasets (with and without social pooling, truncation trick)**

116 ETH-UCY: ETH/UCY dataset consists of 5 scenes eth, hotel, univ, zara1, zara2 extracted from another
 117 source <link> because, in the paper, the source was not mentioned. We Followed the conventional

118 leave-one-out approach, i.e., trained on 4 sets and tested on the last set to get the results. We verified
 119 results within 98% accuracy from claimed results. The dataset was further downsampled by 6 to get
 120 a 0.4 second gap between consecutive frames as demanded by the paper. The result is shown below
 121 in the table. With these 2 experiments, the reduction in error with respect to the previous results
 122 by using CVAE and subsequent reduction by using social pooling layer and truncation trick can be
 123 demonstrated.

Datasets	O-S-TT		PECNet-Ours	
	ADE	FDE	ADE	FDE
ETH	0.58/.57	0.96/.98	0.54/.53	0.87/.87
HOTEL	0.19/.20	0.34/.35	0.18/0.18	0.24/0.23
UNIV	0.39/0.32	0.67/0.53	0.35/0.32	0.60/0.49
ZARA1	0.23/0.23	0.39/0.37	0.22/0.23	0.39/0.35
ZARA2	0.24/0.20	0.35/0.33	0.17/0.20	0.30/0.32

Table 3: Quantitative results obtained versus those of the authors’ (in the form of ours/authors’). ‘Our-S-TT’ represents ablation of our method without social pooling truncation trick. The format for each cell is <claimed result> / <reproduced result>

124 4.3 Change in the structure of CVAE

125 In this experiment during training, the ground truth Eend (G_k) was used to predict the future T_f
 126 instead of the one obtained from the latent variable. We did it on the Stanford drone dataset with
 127 social pooling and got results within 95% accuracy from the claim results.

128 ADE : 10.87 / 10.945

129 FDE : 17.03 / 16.277

130 4.4 Effect of Number of samples (K)

131 We did this experiment on the Stanford drone dataset with social pooling. We trained the PECNet
 132 model with default sigma values and test on different k-sample value with and without truncation.
 133 Without truncation for k-sample \leq 3 we used σ with variance 1 and for k-sample $>$ 3 we used σ with
 134 variance 1.3. With truncation for k-sample $>$ 3 we used σ with variance 1 and for k-sample \leq 3 we
 135 used σ with variance $c * \sqrt{k - 1}$. In this experiment we got results within 95 accuracy from the claim
 136 results.

	1	2	3	5	10	50	100	1000	10000
ADE	24.29	18.457	16.25	14.16	12.04	8.99	8.208	6.81	6.27
FDE	51.84	37.65	32.15	26.73	21.10	12.27	9.73	4.66	2.46
Truncated-ADE	17.62	16.67	15.71	14.788	12.10	8.54	7.70	6.39	6.02
Truncated-FDE	35.02	32.67	30.34	28.57	21.49	11.27	8.54	3.54	1.66

Table 4: Effect of no of samples (K) on ADE, FDE, Truncated-ADE, Truncated-FDE

137 4.5 Conditioned Way-point positions Oracles

138 In this experiment, we conditioned on future trajectory points other than the last observed point,
 139 which we refer to as way-points. This was not clear in the paper about how to calculate FDE error
 140 because we can not predict last observed point in the model so we calculated FDE from the 11th
 141 point of the predicted trajectory. It was done in two parts.

- 142 1. **With oracle:** During prediction of future trajectory (at time of testing and validation), we
 143 gave ground-truth value of conditioned point instead of the best guessed one from sampling
 144 to predict trajectory from the model. The Stanford drone data set with social pooling and
 145 truncation trick was used to match with the results on paper.

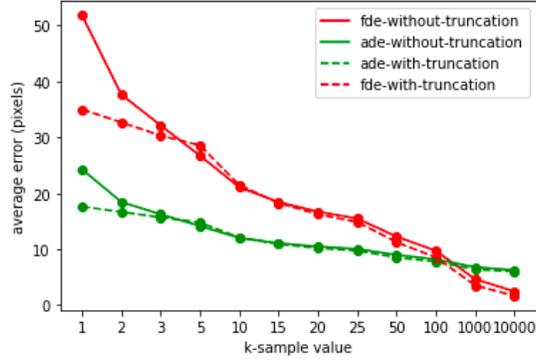


Figure 2: Graph of errors

146
147
148
149

2. **Without oracle:** The same thing was done here except during prediction of the future trajectory the best guess for the conditioned point(predicted by model) was taken (at time of testing and validation). Way-point Prediction Error was calculated as difference between ground truth of conditioned point and the one predicted by the model.

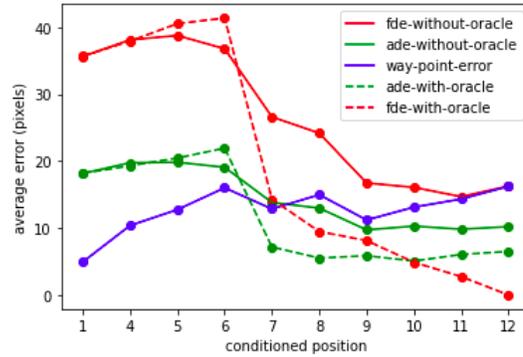


Figure 3: Graph of errors

	1	4	5	6	7	8	9	10	11	12
ADE	18.16	19.76	19.83	19.08	13.82	12.98	9.73	10.29	9.83	10.218
FDE	35.64	38.125	38.77	36.79	26.61	24.18	16.73	16.08	14.69	16.27
Way-point error	4.93	10.38	12.75	16.01	12.86	14.98	11.207	13.12	14.336	16.23
Oracle ADE	18.17	19.30	20.46	21.94	7.17	5.52	5.87	5.074	6.0552	6.51
Oracle FDE	35.68	37.93	40.54	41.38	14.30	9.48	8.13	4.892	2.745	0.0

Table 5: Conditioned Way-point positions and Oracles

150 5 Discussion

151 From each of the experiments, the claims made by the paper as described above can be strongly
152 supported and empirically proved. In order to further study the choice of structure of the network, we
153 performed the following experiments:-

154 5.1 Reference shift <link>[1]

155 We took the reference of the trajectory for each pedestrian as the current point instead of the first point
156 of the past trajectory. This helped the CVAE network to get a better representation of the destination

157 point as all input past trajectories have a common last point, which makes it easier for the encoder and
158 decoder network to function; also, the predictor and social pooling network gets more easily trained.
159 This showed about 8% further decrease in ADE and FDE metrics for drone dataset as follows:-
160 ADE : 8.64
161 FDE : 14.64

162 **5.2 Using encoder and decoder LSTM network instead of MLP <link>[1]**

163 We used encoder LSTM instead of MLP to form the encoding of the past trajectory to accommodate
164 variable length of past trajectory and form a better representation as to the input temporal data.
165 Also, we used the decoder LSTM network to predict the rest of the trajectory given the destination.
166 However, the FDE error reduced by about 5 %, but the ADE is surprisingly more, demonstrating that
167 decoder LSTM does not perform well given the destination point.

168 ADE : 26.9

169 FDE : 14.3

170 **References**

- 171 [1] Lukas Biewald. *Experiment Tracking with Weights and Biases*. Software available from
172 wandb.com. 2020. URL: <https://www.wandb.com/%5C%7D>.
- 173 [2] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. “Crowds by Example”. In: *Computer*
174 *Graphics Forum* 26.3 (2007), pp. 655–664. DOI: [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-8659.2007.01089.x)
175 [8659.2007.01089.x](https://doi.org/10.1111/j.1467-8659.2007.01089.x). eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.1111/](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-8659.2007.01089.x)
176 [j.1467-8659.2007.01089.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-8659.2007.01089.x). URL: [https://onlinelibrary.wiley.com/doi/abs/10.](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-8659.2007.01089.x)
177 [1111/j.1467-8659.2007.01089.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-8659.2007.01089.x).
- 178 [3] Karttikeya Mangalam et al. “It is Not the Journey but the Destination: Endpoint Conditioned
179 Trajectory Prediction”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*.
180 Aug. 2020.
- 181 [4] S. Pellegrini et al. “You’ll never walk alone: Modeling social behavior for multi-target tracking”.
182 In: *2009 IEEE 12th International Conference on Computer Vision*. 2009, pp. 261–268. DOI:
183 [10.1109/ICCV.2009.5459260](https://doi.org/10.1109/ICCV.2009.5459260).
- 184 [5] A. Robicquet et al. “Stanford Drone Dataset”. In: (). URL: [http://cvgl.stanford.edu/](http://cvgl.stanford.edu/projects/uav_data/)
185 [projects/uav_data/](http://cvgl.stanford.edu/projects/uav_data/).