

## Interpretable Social Anchors for Human Trajectory Forecasting in Crowds

Parth Kothari, Brian Siffringer, Alexandre Alahi  
EPFL VITA lab  
CH-1015 Lausanne  
parth.kothari@epfl.ch

### Abstract

Human trajectory forecasting in crowds, at its core, is a sequence prediction problem with specific challenges of capturing inter-sequence dependencies (social interactions) and consequently predicting socially-compliant multimodal distributions. In recent years, neural network-based methods have been shown to outperform hand-crafted methods on distance-based metrics. However, these data-driven methods still suffer from one crucial limitation: lack of interpretability. To overcome this limitation, we leverage the power of discrete choice models to learn interpretable rule-based intents, and subsequently utilize the expressibility of neural networks to model scene-specific residuals. Extensive experimentation on the interaction-centric benchmark TrajNet++ demonstrates the effectiveness of our proposed architecture in expanding its predictions without compromising the accuracy.

### 1. Introduction

Humans naturally navigate through crowds by following the unspoken rules of social motion such as avoiding collisions or yielding right-of-way. Forecasting human motion in public places is a challenging, yet crucial task for the success of many application like deployment of autonomous navigation systems [1, 2, 3], infrastructure planning [25, 21] and evacuation analysis [24, 65]. Therefore, in the last few decades, developing models that can understand human social interactions and forecast their future motions have been an active and challenging area of research.

Early works designed hand-crafted methods based on domain knowledge to forecast human trajectories, either with physics-based models such as Social Forces [23], or with pattern-based models such as discrete choice model (DCM) [7, 26, 31]. These models, based on domain knowledge, were successful in showcasing crowd phenomena like collision avoidance and leader-follower type behaviors.

To appear in Computer Vision and Pattern Recognition (CVPR) 2021

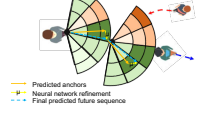


Figure 1: While navigating in crowds, humans display various social phenomena like collision avoidance (from red trajectory) and leader-follower (towards blue trajectory). We present a model that not only outputs accurate future trajectories but also provides a high-level rationale behind its predictions, owing to the interpretability of discrete choice models. (Un)favorable anchors are shown in green (red).

Moreover, the hand-designed nature of these models rendered their predictions to be interpretable. However, human motion in crowds is much more complex and due to its long-term nature, these first-order models suffer from predicting inaccurate trajectories. Building on the success of recent neural network-based methods in learning complex functions and long-term dependencies, Alahi et al. [4] proposed the first neural network (NN) based trajectory forecasting model, Social LSTM, which outperformed the hand-crafted methods on distance-based metrics. Due to the success of Social LSTM, neural networks have become the de-facto choice for designing human trajectory models [2, 35, 65, 28, 19]. However, current NN-based trajectory forecasting models suffer from a significant limitation: lack of interpretability regarding the model’s decision-making process. In this work, we are interested in combining the best of two paradigms of human trajectory forecasting (see

Fig. 1): the interpretability of the trajectories predicted by hand-crafted models, in particular discrete choice models [7, 35], and the high accuracy of the neural networks based predictions. With this objective, we propose a model that provides a probability distribution over a discrete set of possible future intents. This is designed as a function of the pedestrian’s speed and direction of movement. Our model learns the probability distribution over these intents with the help of a choice model augmented owing to its ability to output interpretable decisions. To this end, we resort to a novel hybrid and interpretable framework in DCM [35], where knowledge-based hand-crafted functions can be augmented with neural network representations, without compromising the interpretability.

Our architecture augments each predicted high-level intent with a scene-specific residual term generated by a neural network. The advantage of this two-fold: first, the residual allows to expand the output space of the model from a discrete distribution to a continuous one. Secondly, it helps to incorporate the complex social interactions as well as the long-term dependencies that the first-order hand-crafted models fail to capture, leading to an increase in prediction accuracy. Overall, we can view our architecture as disentangling high-level coarse intents and lower-level scene-specific nuances of human motion.

We demonstrate the efficacy of our proposed architecture on TrajNet++ [23], an interaction-centric human trajectory forecasting benchmark comprising of well-sampled real-world trajectories that undergo various social phenomena. Through extensive experimentation, we demonstrate that NN-based models outperform competing baselines on both real-world and synthetic datasets, while at the same time providing a rationale behind high-level decisions, an essential component required for safety-critical applications like autonomous systems.

### 2. Related Work

#### 2.1. Social Interactions

Current human trajectory forecasting research can be categorized into learning human motion (social) interactions and human-space (physical) interactions. In this work, we focus on the task of discrete choice models that aim at understanding social interactions in crowds. The human social interactions are usually modeled either using knowledge-based models or using neural networks. Knowledge-based Models: With a specific focus on pedestrian path forecasting problems, Helbing and Molnar [23] presented a force-based motion model with attractive forces (towards the goal and one’s own group) and repulsive forces (away from obstacles), called Social Force model. Bursztke et al. [14] utilize the cellular automaton model to predict pedestrian motion by dividing the environment into

uniform grids and assigning transition probability matrices to the pedestrians. Similarly, discrete choice modelling yields a grid for selecting the next action, but relative to each individual [7, 26, 31]. The high interpretability and design flexibility of DCM allowed its application to many topics such as pedestrian flows [49], walking in groups [46, 63], collision avoidance [8, 40], and critical or emergency situations [24, 45, 61]. Human social interactions have also been modelled using other knowledge-based perspectives [57, 49, 9]. While the hand-crafted functions of these methods lead to interpretable outputs, they are often too simple to capture the complexity of human interactions.

Neural Network-based Models: In the past few years, methods based on neural networks (NNs) that infer social interactions in a data-driven fashion have been shown to outperform the knowledge-based works on distance-based metrics. Social LSTM [4] introduced a novel social pooling layer to capture social interactions of nearby pedestrians. Various other interaction-capturing NN modules have been proposed in literature [45, 15, 12, 21, 44, 56, 38, 56]. To provide different weights to neighbors that affect the trajectory of the pedestrian of interest, multiple weights [62, 18, 55, 53, 31, 6, 17, 27, 44, 50] propose to utilize attention mechanisms [58, 9]. The attention weights are either learned or handcrafted based on domain knowledge (e.g., euclidean distance). However, these data-driven methods lack the ability to output predictions that can be explained, unlike their knowledge-based counterparts.

In this work, we combine the strengths of rule-based models to output high-level intents that are interpretable, and NN-based models to capture complex social interactions that take into account the long-term motion characteristics.

#### 2.2. Multimodality

Training neural networks based on minimization of  $L_2$  loss leads to the model outputting the mean of all the outcomes. One solution to multimodal forecasting is to explicitly output multiple modes using the decoder architecture, for instance, using Mixture Density Networks [13]. However, this training technique suffers from numerical instabilities, often leading to mode collapse.

Another recent popular approach is based on generative modelling [3, 28, 41, 7]. Generalized autoregressive denoising (GAD) framework [10] uses a generative model to learn the probability distribution of the future trajectories conditional to the past scene, thereby naturally offering a possibility to output multiple samples. Lee et al. [33] propose a recurrent encoder-decoder architecture for pedestrian trajectory forecasting. Helbing and Molnar [23] present a force-based motion model with attractive forces (towards the goal and one’s own group) and repulsive forces (away from obstacles), called Social Force model. Bursztke et al. [14] utilize the cellular automaton model to predict pedestrian motion by dividing the environment into

Model	ADE/FDE	Col	Top-1 ADE/FDE
S-LSTM [4]	0.291/0.24	5.3	0.373/0.28
WPA [16]	0.280/0.24	4.8	0.329/0.24
SGAN [1]	0.278/0.24	5.1	0.324/0.24
CVAN [18]	0.280/0.24	5.2	0.329/0.24
MMK [6]	0.287/0.24	5.2	0.324/0.24
Sancher (Ours)	0.278/0.24	4.8	0.329/0.24

Table 1: Performance on TrajNet++ synthetic data. Errors reported are ADE / FDE in meters. Col is % in %. We observe the trajectories for 9 times-steps (3.6 secs) and performance prediction for the next 12 (4.8 secs) time-steps. \*Unimodal.

Model	ADE/FDE	Col	Top-1 ADE/FDE
S-LSTM [4]	0.291/0.24	5.3	0.373/0.28
WPA [16]	0.280/0.24	4.8	0.329/0.24
SGAN [1]	0.278/0.24	5.1	0.324/0.24
CVAN [18]	0.280/0.24	5.2	0.329/0.24
MMK [6]	0.287/0.24	5.2	0.324/0.24
Sancher (Ours)	0.278/0.24	4.8	0.329/0.24

Table 2: Performance on TrajNet++ real data. Errors reported are ADE / FDE in meters. Col is % in %. We observe the trajectories for 9 times-steps (3.6 secs) and performance prediction for the next 12 (4.8 secs) time-steps. \*Unimodal.

We demonstrate the ability of our network to output interpretable intents in Fig. 3. The direction of the pedestrian’s intent is normalized and is facing towards the right. For each row in Fig. 3, in addition to the ground-truth map (left-most), we illustrate the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene. In the first row, we observe that the model correctly chooses an intent, while the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene. In the first row, we observe that the model correctly chooses an intent, while the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene. In the first row, we observe that the model correctly chooses an intent, while the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene.

#### Evaluation: we consider the following metrics:

- Average Displacement Error (ADE):** the average  $L_2$  distance between ground-truth and model prediction over predicted time steps.
- Final Displacement Error (FDE):** the distance between the final predicted destination and the ground-truth destination at the end of the prediction period.
- Collision 1 - Prediction collision (Col) [12]:** this metric calculates the percentage of collision between the pedestrian of interest and the neighbors in the predicted scene. This metric indicates whether the predicted model trajectories collide, i.e., whether the model learns the notion of collision avoidance.
- Top-3 ADE/FDE:** given  $t$  output predictions for an observed scene, this metric calculates the ADE/FDE of the prediction closest to the ground-truth trajectory in terms of ADE.

Baselines: we compare against the following baselines: 1. **S-LSTM**: we compare to S-LSTM [4] baseline that uses a unimodal trajectory distribution.

2. **Winner-Takes-All (WTA)**: this architecture was proposed in [51] to encourage the network to output discrete trajectories of all the pedestrians.

3. **SGAN**: Social GAN [21], a popular generative model to tackle multimodal trajectory forecasting of the intent is normalized and is facing towards the right. For each row in Fig. 3, in addition to the ground-truth map (left-most), we illustrate the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene. In the first row, we observe that the model correctly chooses an intent, while the activation maps of all combined factors, the neural network (NN) map, the overall DCM map and finally the dominant baseline rules that comprise the DCM function, according to the presented scene.

4. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

5. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

6. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

7. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

8. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

9. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

10. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

11. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

12. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

13. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

14. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

15. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

16. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

17. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

18. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

19. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

20. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

21. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

22. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

23. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

24. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

25. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

26. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

27. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

28. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

29. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

30. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

31. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

32. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

33. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

34. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

35. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

36. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

37. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

38. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

39. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

40. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

41. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

42. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

43. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

44. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

45. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

46. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

47. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

48. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

49. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

50. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

51. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

52. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

53. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

54. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

55. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

56. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

57. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

58. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

59. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

60. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

61. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

62. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

63. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

64. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

65. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

66. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

67. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

68. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

69. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

70. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

71. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

72. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

73. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

74. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

75. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

76. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

77. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

78. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

79. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

80. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

81. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

82. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

83. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

84. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

85. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

86. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

87. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

88. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

89. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

90. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

91. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

92. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

93. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

94. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

95. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

96. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

97. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

98. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

99. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

100. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

101. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

102. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

103. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

104. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

105. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

106. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (discussed next) to output high-level intents.

107. **Neural Network (NN)**: we compare against the NN baseline that directly outputs the NN residuals without any prior anchors.

108. **Sancher (Ours)**: our proposed method that utilizes 15 scene-specific residuals and 15 scene-specific residuals (