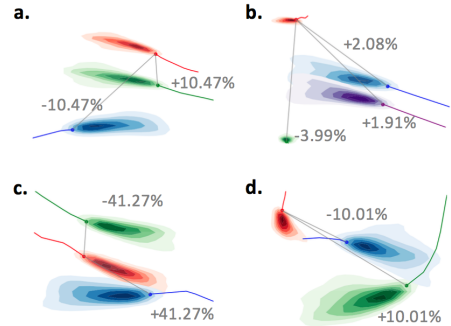


1 **R1, R2: Technical novelty.** While our work builds on top of previously proposed methods including GAN, attention
 2 mechanisms, and GAT, their adaptation, integration and application to the complex problem of human trajectory
 3 prediction is a non-trivial technical challenge. Our main technical contribution is a novel neural model that 1) better
 4 encodes social cues (a crucial factor to predict human trajectories) based on neural graph structures, and 2) enables
 5 multimodal predictions (an intrinsic property of future human motion) inspired by BicycleGAN architecture. Learning
 6 an optimal representation of social behavior is a non-trivial task, as humans follow unwritten social rules and juggle a
 7 variety of implicit factors. As we discuss in our literature review and experimental evaluation, prior works that model
 8 social interactions fall short by either limiting the expressiveness of the models (eg: by using a pooling mechanism that
 9 is unable to capture interactions in large scenes), or by imposing human-defined constraints rather than learning from
 10 the data. Further, applying a multi-modal distribution to represent future human trajectories is not easy since it requires
 11 solving a delicate trade-off between increasing variance of each model vs. increasing multimodality, which relates to
 12 the degradation of the model’s predictions as a function of allowed samples.

13 **R1, R3: Analysis on the attention weights.** As the reviewers suggested, we
 14 have run several experiments exploring the correlation between the weights
 15 and different pedestrian features. The figure on the right depicts the attention
 16 weights of other pedestrians wrt. the red pedestrian for different scenes.
 17 Weight values correspond to the percent decrease/increase compared to without
 18 attention: positive indicates more attention was paid to that interaction, and
 19 negative indicates less attention was paid. From Scenes (a, b) we can infer
 20 that Euclidean distance is one feature the network implicitly uses to assign
 21 attention. Scenes (c) and (d) show however that the attention further generalizes
 22 in learning which agents are important socially: in Scene (c) it pays large
 23 attention to the blue agent it may collide with in the future, even though that
 24 agent is farther away from it than the green one, and in Scene (d) it ignores
 25 the blue agent for the farther green agent with whom it might collide with. The diversity of social awareness that the
 26 model displays validates the choice of GAT over prior works. We will include this figure and a visualization of physical
 27 attention to the final manuscript, with more samples in the supplementary material.



28 **R1: Improved generalization claim.** We thank R1 for their suggestion and have updated the table as shown below. We
 29 see that the GAT architecture initially performs better than the BiGAN, but with fewer samples the GAT error increases
 30 faster. This aligns with our intuition that the inclusion of the BiGAN in our architecture enables for better capturing of
 31 a multimodal distribution, instead of generating samples from a unimodal distribution. We note that BiGAN’s error
 32 increases slightly slower than Social-BiGAT, but the inclusion of the GAT allows for a lower overall error.

33 **R1: Comparison with SOTA baselines.**

Model	K=20	K=10	K=5	K=1	% Increase
GAT	0.518 / 1.064	0.529 / 1.127	0.584 / 1.241	0.682 / 1.494	31.6% / 40.4%
BiGAN	0.523 / 1.091	0.531 / 1.144	0.579 / 1.298	0.662 / 1.439	26.6% / 31.9%
Social-BiGAT	0.476 / 0.998	0.488 / 1.096	0.527 / 1.260	0.606 / 1.328	27.3% / 33.1%

34 Yes, Sophie was the state of the art in terms
 35 of ADE/FDE at the time of submission.

36 **R1: Clarifying pooling for joint feature
 37 representation.**

In architectures that use pooling, each pedestrian is encoded into a single vector. Across all pedestrians
 38 the vectors are pooled, resulting in a single vector representing the social interactions between all agents in the scene
 39 (not specific to any particular agent), which each pedestrian receives. Our choice of GAT serves as a major improvement
 40 as it allows for different interactions to be attended to based on their social importance, which the network learns.

41 **R2: Inputting top-view image instead of image sequences.** Thanks to the reviewer for pointing out the confusion,
 42 we will clarify this in the paper. However, we would like to point out that most prior research/datasets in this field are
 43 using a single top-down image input, such as CAR-Net, Desire, and Sophie, and the more recent Social Ways, and
 44 Multi Agent Tensor Flow. Moreover, our approach has also been tested on sequences where the scene’s angle view is
 45 not necessarily perfectly top-view, e.g. UCY dataset. The criticism of not accepting a sequences of images is valid, as
 46 we only use a single image. This is done for fair comparison, as prior methods also only use single images. Finally,
 47 our work is indeed limited by requiring temporal sequences of pedestrian locations. However, while we agree that
 48 performing end-to-end tracking and forecasting directly on sequences of images is a very good future direction, we are
 49 currently limited by a lack of proper benchmark data. All existing benchmarks follow the same input modality, and we
 50 do the same to make the comparison feasible. In the future we definitely hope to research this same problem end-to-end.

51 **R2, R3: Lack of clarity in notations, architecture.** We will clarify the equations and better explain the variables we
 52 introduce before using them, along with labelling weights in Fig. 2, which should improve the architecture explanation.
 53 We also will add an in-depth ablative analysis and explanation of model components in the supplementary material.

54 **R2: Better visualizations & demos.** As suggested by R2, we have already started implementing it by adding several
 55 examples to the supplementary work, where we can visualize multiple dimensions of the latent space across multiple
 56 pages and better illustrate how our model qualitatively results in multimodality. We will also add figures exploring
 57 how the latent noise generated through the encoder changes when varying parts of the scene (such as number of agents,
 58 or speed and direction of agents). We also are looking forward to publishing demo videos that we can place in the
 59 supplementary material as well. Space permitting, we will add these figures to the main paper as well.