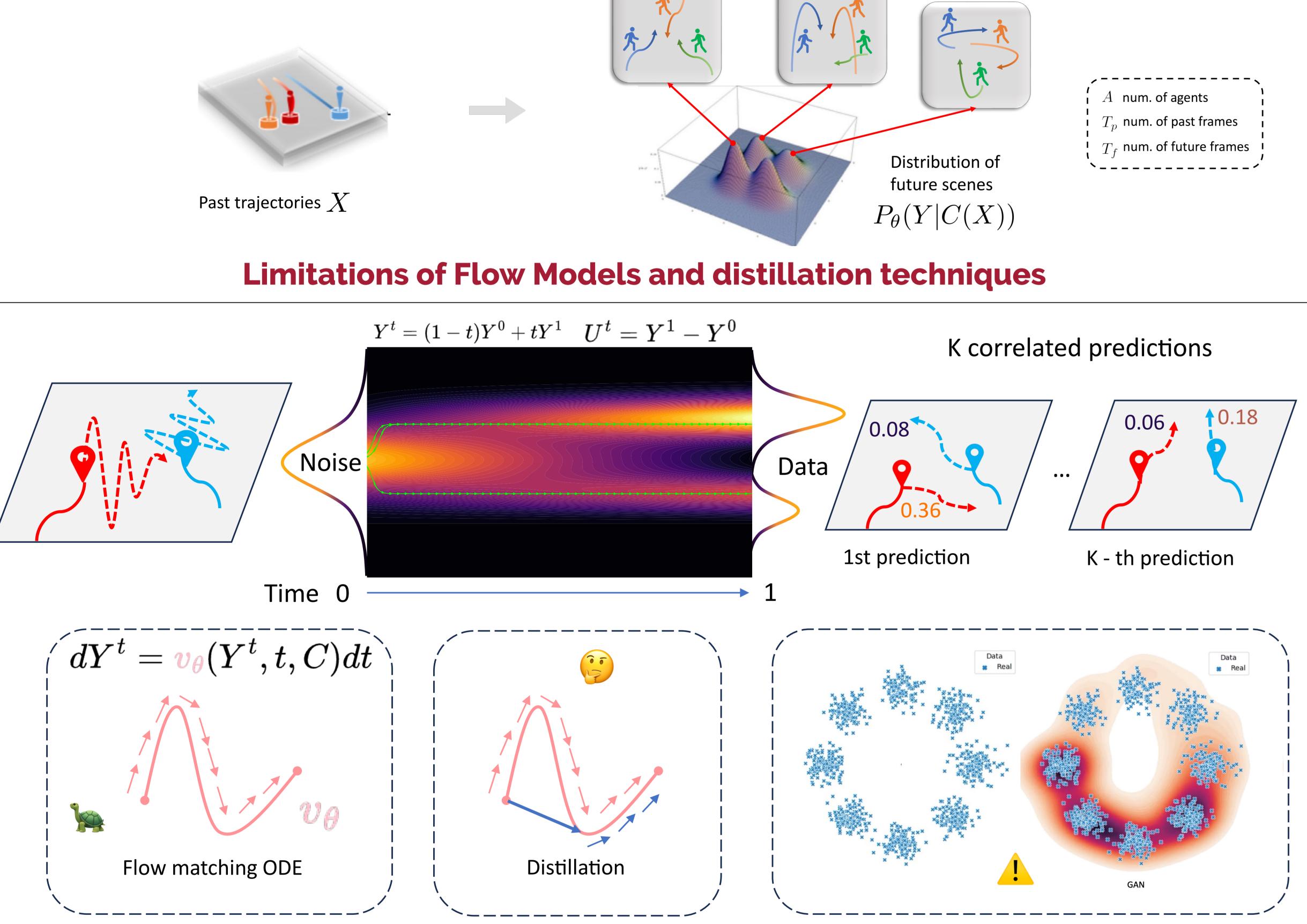


# MoFlow: One-Step Flow Matching for Human Trajectory Forecasting via Implicit Maximum Likelihood Estimation based Distillation

### Background

**Problem formulation:** Given the past trajectory history  $X = [\mathbf{x}_1^\top; \mathbf{x}_2^\top; \dots; \mathbf{x}_A^\top] \in \mathbb{R}^{A \times 2T_p}$  of the agents, we want to model the distribution of the future scenes  $Y \sim P_{\theta}, Y \in \mathbb{R}^{A imes 2T_f}$  jointly so that the generated samples are socially and physically compliant. Note that A represents the number of agents and  $T_p, T_f$  are the number of past and future frames respectively.



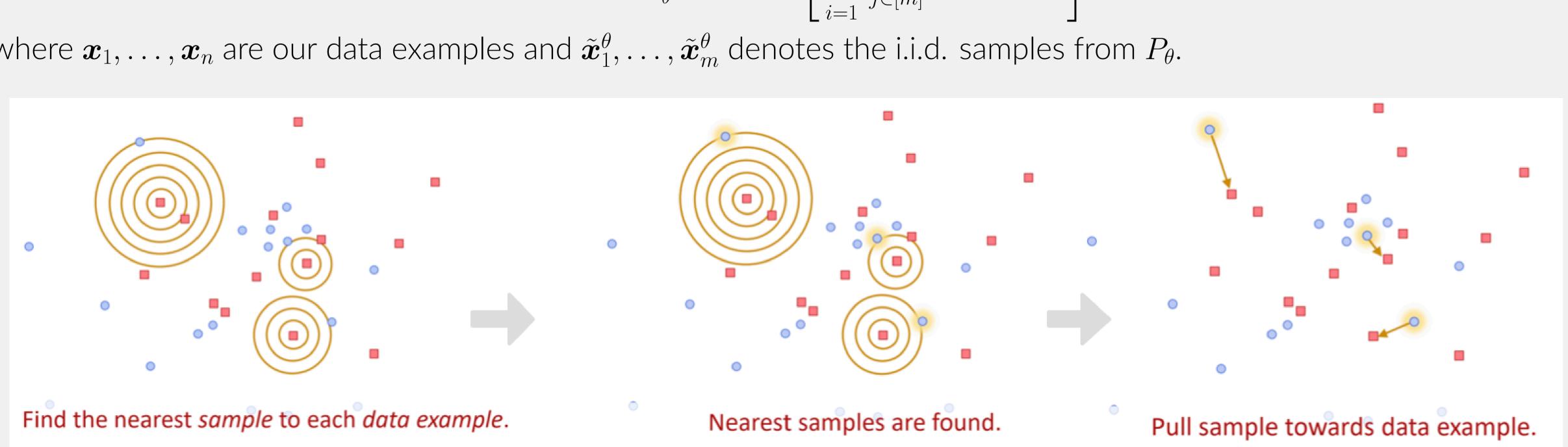
- The numerical ODE solver requires numerous iterative network evaluations, making the inference **slow** and **expensive**.
- Most distillation methods fail to preserve good quality of samples from teacher model.
- Consistency Distillation requires samples  $Y^{t_n}$  from the teacher model at different noise levels for all timestamps.

### **Implicit Maximum Likelihood Estimator**

The Implicit Maximum Likelihood Estimator (IMLE) is defined as

$$\mathsf{MLE} := \arg\min_{\theta} \mathbb{E}_{\tilde{\boldsymbol{x}}_{1}^{\theta}, \dots, \tilde{\boldsymbol{x}}_{m}^{\theta}} \left[ \sum_{i=1}^{n} \min_{j \in [m]} \| \tilde{\boldsymbol{x}}_{j}^{\theta} - \boldsymbol{x}_{i} \| \right]$$

where  $m{x}_1,\ldots,m{x}_n$  are our data examples and  $ilde{m{x}}_1^ heta,\ldots, ilde{m{x}}_m^ heta$  denotes the i.i.d. samples from  $P_ heta$ .



<sup>1</sup>University of British Columbia

## Contributions

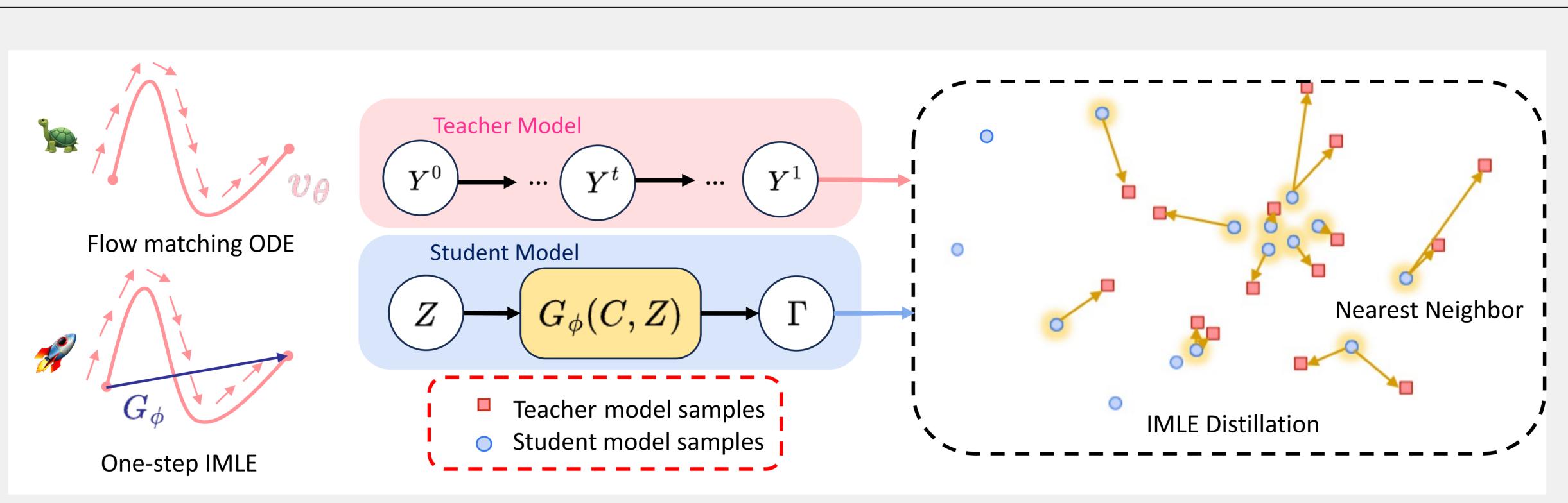
- 1. We present a novel **Mo**tion prediction **Flow** matching model for trajectory prediction tasks. Also, we design a novel human trajectories.
- 2. We propose a **one-step** distillation method for flow models based on IMLE. Our approach accelerates the standard conditional flow matching inference by **100x** without sacrificing quality of generated trajectories.
- 3. Both MoFlow and the proposed distillation method deliver SOTA performance over three human motion datasets, producing diverse trajectories that are physically and socially plausible.

### **MoFlow Objective**

Our model outputs K scene-level waypoint predictions, denoted by  $\{S_i\}_{i=1}^K, S_i \in \mathbb{R}^{A \times 2T_f}$ , alongside corresponding **classification logits**  $\{\zeta_i\}_{i=1}^K, \zeta_i \in \mathbb{R}$ . Building upon original Flow Matching loss, we derive  $\bar{\mathcal{L}}_{\mathsf{FM}} = \mathbb{E}_{Y^{t}, Y^{1}, t} \left[ \|S_{j^{*}} - Y^{1}\|_{2}^{2} + \operatorname{CE}(\zeta_{1:K}, j^{*}) \right], \quad j^{*} = \arg\min_{j} \|S_{j} - Y^{1}\|_{2}^{2},$ 

where  $CE(\cdot, \cdot)$  denotes the cross-entropy loss.

### **IMLE** Distillation



We adapt IMLE to a new objective tailored for our conditional generation task, expressed as follows

$$\hat{\phi}_{\mathsf{IMLE}} := \arg\min_{\phi} \mathbb{E}_{Z_1, \dots, Z_m} \left[ \min_{j \in [m]} \mathcal{L}_{\mathsf{IMLE}}(\hat{Y}_{1:K}^1, G_{\phi}(C, Z_j)) \right], \\ \mathcal{L}_{\mathsf{IMLE}}(\hat{Y}_{1:K}^1, \Gamma) = \frac{1}{K} \left( \sum_{i=1}^K \min_{j} \| \hat{Y}_i^1 - \Gamma^{(j)} \| + \sum_{j=1}^K \min_{i} \| \hat{Y}_i^1 - \Gamma^{(j)} \| \right),$$

where our conditional IMLE generator  $G_{\phi}$  uses a noise vector among  $Z_1, \ldots, Z_m \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and context C to generate K-component trajectories  $\Gamma \in \mathbb{R}^{K \times A \times 2T_f}$ , matching the shape of  $\hat{Y}_{1:K}^1$  from our teacher model.

### **Quantitative Results on NBA/ETH-UCY/SDD Datasets**

Time	MemoNet [54]	NPSN [4]	GroupNet [53]	MID [16]	LED <sup>†</sup> [33]	LED* [33]	MoFlow	IMLE
1.0s	0.38/0.56	0.35/0.58	0.26/0.34	0.28/0.37	<u>0.21</u> /0.28	<b>0.18</b> / <u>0.27</u>	0.18/0.25	0.18/0.25
2.0s	0.71/1.14	0.68/1.23	0.49/0.70	0.51/0.72	0.44/0.64	0.37/ <u>0.56</u>	0.34/0.47	<u>0.35</u> / <b>0.47</b>
3.0s	1.00/1.57	1.01/1.76	0.73/1.02	0.71/0.98	0.69/0.95	<u>0.58/0.84</u>	0.52/0.67	0.52/0.67
Total (4.0s)	1.25/1.47	1.31/1.79	0.96/1.30	0.96/1.27	0.94/1.21	0.81/1.16	0.71/0.87	0.71/0.87

Subsets	MID [16]	GroupNet [53]	TUTR [42]	EqMotion [55]	EigenTraj [ <mark>5</mark> ]	LED* [33]	SingularTraj [ <mark>6</mark> ]	MoFlow	IMLE
ETH	0.39/0.66	0.46/0.73	0.40/0.61	0.40/0.61	0.36/0.53	0.39/0.58	0.35/0.42	0.40/0.57	0.40/0.58
HOTEL	0.13/0.22	0.15/0.25	<b>0.11</b> / <u>0.18</u>	<u>0.12/0.18</u>	<u>0.12</u> /0.19	0.11/0.17	0.13/0.19	0.11/0.17	0.12/0.18
UNIV	<b>0.22</b> /0.45	0.26/0.49	<u>0.23</u> /0.42	<u>0.23</u> /0.43	0.24/0.34	0.26/0.44	0.25/0.44	<u>0.23</u> / <b>0.39</b>	<u>0.23</u> / <b>0.39</b>
ZARA1	0.17/0.30	0.21/0.39	0.18/0.34	0.18/0.32	0.19/0.33	0.18/0.26	0.19/0.32	0.15/0.26	<u>0.16</u> / <b>0.26</b>
ZARA2	<u>0.13</u> /0.27	0.17/0.33	<u>0.13</u> /0.25	<u>0.13/0.23</u>	0.14/0.24	<u>0.13</u> / <b>0.22</b>	0.15/0.25	0.12/0.22	<u>0.13</u> / <b>0.22</b>
AVG	<u>0.21</u> /0.38	0.25/0.44	<u>0.21</u> /0.36	0.32/0.35	<u>0.21</u> /0.34	<u>0.21/0.33</u>	<u>0.21</u> / <b>0.32</b>	0.20/0.32	<u>0.21/0.33</u>

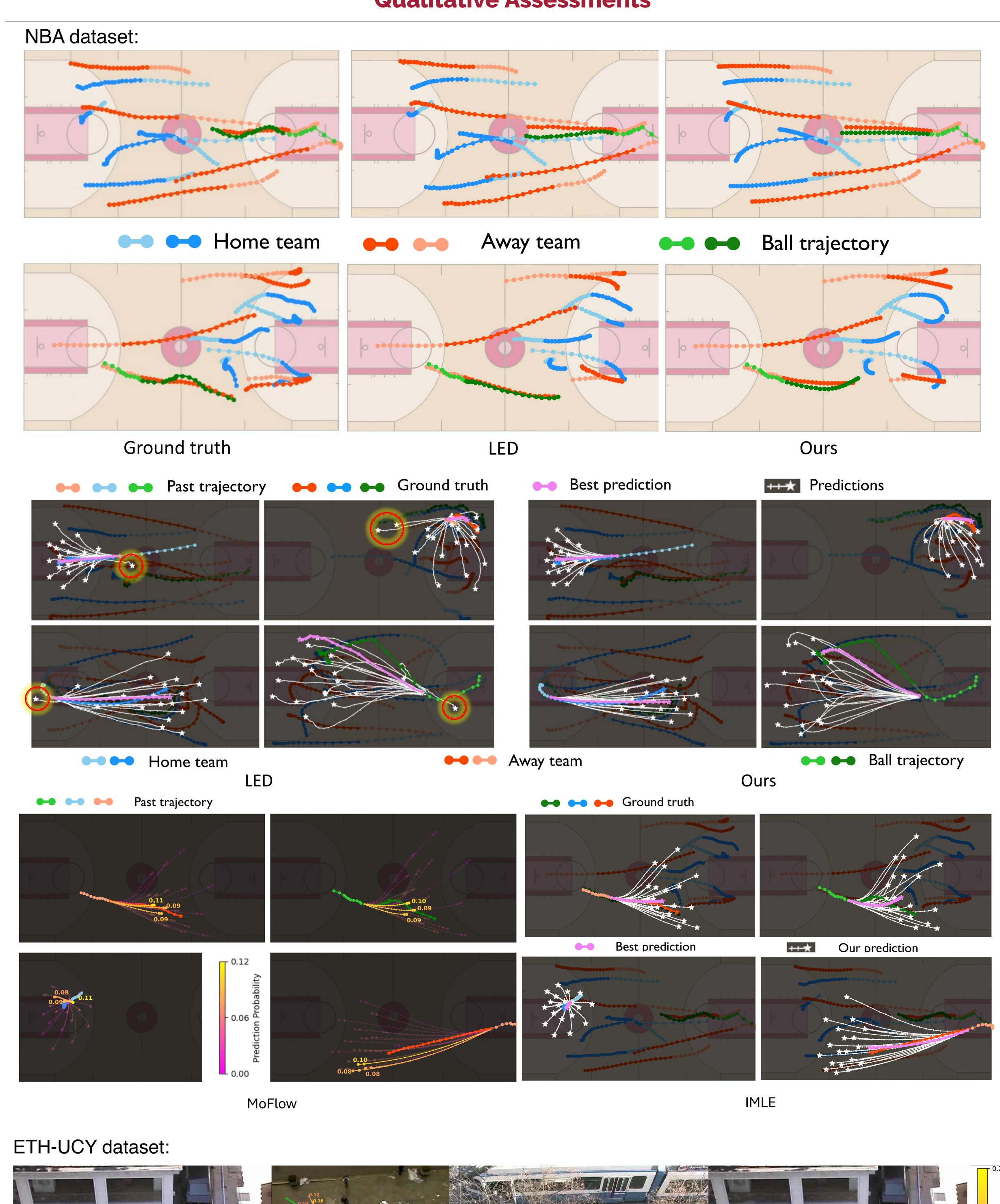
Yuxiang Fu<sup>1,2</sup> Qi Yan<sup>1,2</sup> Lele Wang<sup>1</sup> Ke Li<sup>4</sup> Renjie Liao<sup>1,2,3</sup>

<sup>2</sup>Vector Institute for Al <sup>3</sup>Canada CIFAR AI Chair <sup>4</sup>Simon Fraser University

flow matching loss that encourages learning a diverse set of future trajectories that well capture the multi-modality of

Method	ADE/FDE
Evo-Graph [23]	13.90/22.90
Y-net [32]	11.49/20.23
GroupNet [53]	9.31/16.11
CAGN [10]	9.42/15.93
NPSN [4]	8.56/ <u>11.85</u>
MemoNet [54]	9.50/14.78
SocialVAE [57]	8.88/14.81
MID [16]	9.73/15.32
TUTR [42]	<u>7.76</u> /12.69
LED [33]	8.48/ <b>11.66</b>
MoFlow	<b>7.50</b> /11.96
IMLE	7.85/12.86

.









### **Qualitative Assessments**