

## Learning to Simulate Self-driven Particles System with **Coordinated Policy Optimization** Code and demo video decisionforce.github.io/CoPO/

Zhenghao Peng<sup>§</sup>, Quanyi Li<sup>‡</sup>, Chunxiao Liu<sup>†</sup>, Bolei Zhou<sup>§</sup> <sup>§</sup> CUHK <sup>†</sup> SenseTime Research <sup>‡</sup> Centre for Perceptual and Interactive Intelligence, CUHK

## Task

- Realistic Crowd Actions
- Safe Driving
- Social Behaviors

Environments powered by:

**METADRIVE** 



## Method

Step 1: Local Coordination for each policy

Step 2: Global Coordination to update global LCF



Local Coordination: Update policies to maximize coordinated reward Coordinated Reward:  $r_1^C = \cos(\Phi)r_1^I + \sin(\Phi)r_1^N$ , where  $\Phi \in [-90^\circ, 90^\circ]$  is LCF Maximize this reward via PPO loss

**Global Coordination: Adjust LCF to maximize global reward** Meta-gradient to update  $\Phi$ :  $\nabla_{\Phi} J^{G}(\theta^{new}) = \nabla_{\theta^{new}} J^{G}(\theta^{new}) \nabla_{\Phi} \theta^{new}$