# Smart TASEP

Deep Q-Learning in Intelligent Matter Simulations

Jonas Märtens

February 8, 2024

LMU Munich

# Introduction

Smart TASEP: Deep Q-Learning in Intelligent
Matter Simulations

- TASEP: Totally Asymmetric Simple
  Exclusion Process
- Smart TASEP: TASEP with intelligent
  agents
  $\rightarrow$ Intelligent Matter Simulations
- Deep Q-Learning: Reinforcement
  Learning with Neural Networks

cool TASEP teaser pic

Introduction
○●○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○

Conclusion and Outlook
○○

TASEP: Totally Asymmetric Simple Exclusion Process

- Stochastic Process on a grid lattice
- Exclusion: Only one particle per site
- Simple: Jumps only one cell far
- Asymmetric: Only to the right

Possible with either

- Periodic Boundary Conditions
- Insertion and Extraction rates at the boundaries

both TASEP pics

Introduction
○○●○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○○

Conclusion and Outlook
○○

# Modifications

Two modifications (before smart)

## (T)ASEP

- Symmetric in one direction and totally asymmetric in the other
- Directed flow with multiple lanes
- compare traffic flow, kinesin transport
- motivates transport optimization

TASEP pic with speeds, maybe kinesin

## Speeds

- Each particle has a speed
- Speed is probability to jump
- Drawn from a distribution
- If different, jams

Introduction
○○○●

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○○

Conclusion and Outlook
○○

# Make it Smart!

Smart TASEP: TASEP with intelligent agents

- have a goal (e.g. go forward)
- sense the environment
- act according to the goal and the environment
- learn from the environment themselves

$\rightarrow$ Intelligent Matter Simulation
$\rightarrow$ Reinforcement Learning - reward based

Observation distance pic? Leading to action?

Introduction
○○○○

Make it Smart!
○●

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○○

Conclusion and Outlook
○○

# Deep Q-Learning

## Deep Q-Learning Algorithm

- Reinforcement learning with deep neural networks
- Model-free: No knowledge of the environment
- Q-Learning: Value-based method
- Off-policy: Learns from old experiences

DQL components

Introduction
OOOO

Make it Smart!
OO

**Deep Q-Learning**
O●OO

Implementation
OOOOO

Baseline
OOOO

Naive Policy
OO

Smarticles
OOOOOOOOOOOOO

Conclusion and Outlook
OO

large picture of flattened observation fed into a neural network, leading to a Q-value for each action

Introduction
OOOO

Make it Smart!
OO

**Deep Q-Learning**
OO●O

Implementation
OOOOO

Baseline
OOOO

Naive Policy
OO

Smarticles
OOOOOOOOOOOOO

Conclusion and Outlook
OO

## How does the network learn?

- Neural Network $\rightarrow$ backpropagation on loss
- Problem: no target value for Q-Learning
- Solution: Use current best guess as target:
  $y = r + \gamma \max_{a'} \hat{Q}(s', a')$
- Target network with soft updates to stabilize learning
- Use gradients for AdamW optimizer
- Replay buffer
  - Sample efficiency
  - Break correlation

backpropagation pic

Introduction
0000

Make it Smart!
00

Deep Q-Learning
000●

Implementation
00000

Baseline
0000

Naive Policy
00

Smarticles
000000000000

Conclusion and Outlook
00

# Implementation

# The `smarttasep` package

## Features

- Easy installation with `pip`
- Fully Customizable
  - Environment size
  - Reward function
  - Network architecture
  - Hyperparameters
  - Training parameters
  - Algorithm choices
- Real-time visualization
- Interactive evaluation
- Management of experiments

## `smarttasep` components

- `smarttasep.GridEnv`: The 2D (T)ASEP environment
- `smarttasep.DQN`: The neural network
- `torchrl.ReplayBuffer`: Efficient, tensor-based replay buffer
- `smarttasep.Trainer`: Wrapper class for training, simulation and experiment management
- `smarttasep.Playground`: Real-time interactive evaluation of trained agents
- `smarttasep.EnvParams`, `smarttasep.Hyperparams`: Configuration classes

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

**Implementation**
○●○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○

Conclusion and Outlook
○○

Let's have a look at some examples! (show usage videos)

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

**Implementation**
○○●○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○

Conclusion and Outlook
○○

For each experiment, the Hyperparameters should be optimized.

### Parameters

- **Learning rate** of the AdamW optimizer.
- **Discount factor** $\gamma$.
- **Replay buffer size**.
- **Batch size** for training.
- **Target network update rate** $\tau$.
- **Exploration-exploitation tradeoff** via the decay constant $\eta$ of the exploration rate $\epsilon(t) = \epsilon_{\text{end}} + (\epsilon_{\text{start}} - \epsilon_{\text{end}})e^{-n_{\text{steps}}/\eta}$.
- **Neural network architecture** (number of hidden layers and neurons per layer).
- **Activation function** of the neural network.

hidden layers sizes example

Introduction
0000

Make it Smart!
OO

Deep Q-Learning
0000

Implementation
OOOO●

Baseline
0000

Naive Policy
OO

Smarticles
000000000000

Conclusion and Outlook
OO

# Baseline

# Overview of Experiments

## Goals

- Current optimization
- General insights into intelligent matter simulations

## Experiments

1. **Baseline**: Classical 2D TASEP with speeds
2. **Naive current optimization**: Hard-coded agents
3. **Smart TASEP**: Deep Q-Learning agents with slightly different goals
   - **Simple reward**: Go forward
   - **Complex reward**: Go forward while forming lanes

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

**Baseline**
○●○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○

Conclusion and Outlook
○○

## Setup

- Classical 2D TASEP
- Periodic boundary conditions
- Checkerboard configuration

Introduction
0000

Make it Smart!
00

Deep Q-Learning
0000

Implementation
00000

**Baseline**
00●0

Naive Policy
00

Smarticles
000000000000

Conclusion and Outlook
00

## Setup

- Classical 2D TASEP
- Periodic boundary conditions
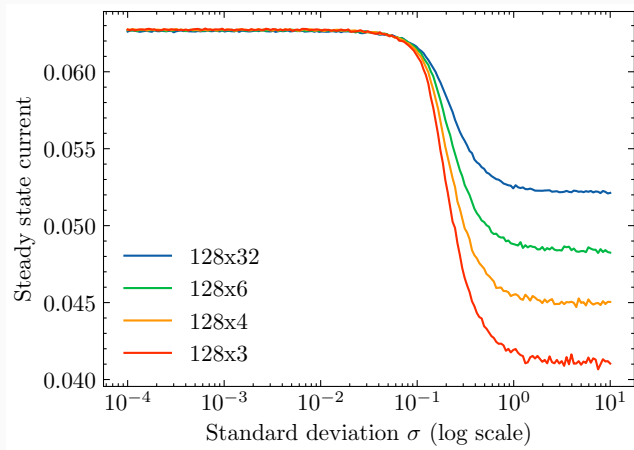- Checkerboard configuration
- Normally distributed speeds



Speeds of particles in the baseline experiment.

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○●○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○○

Conclusion and Outlook
○○

## Setup

- Classical 2D TASEP
- Periodic boundary conditions
- Checkerboard configuration
- Normally distributed speeds
- Waiting for steady state



Equilibration of the current

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○●○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○○○

Conclusion and Outlook
○○

## Setup

- Classical 2D TASEP
- Periodic boundary conditions
- Checkerboard configuration
- Normally distributed speeds
- Waiting for steady state
- Measuring average current



Steady state current as a function of $\sigma$ for different system sizes. 800 runs and 150k steps averaged.

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○●○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○

Conclusion and Outlook
○○

## Conditions for a forward move

1. randomly picked cell is occupied ($p_\text{occ} = \rho$)

2. speed allows move ($p_\text{spd} = \bar{v} = \mu$)

3. forward direction is picked ($p = 0.5$)

4. next cell is empty ($p_\text{emp} = 1 - \rho$)*

$$\implies \langle J \rangle = p \cdot \mu \cdot \rho \cdot (1 - \rho)$$
$$= 0.0625$$



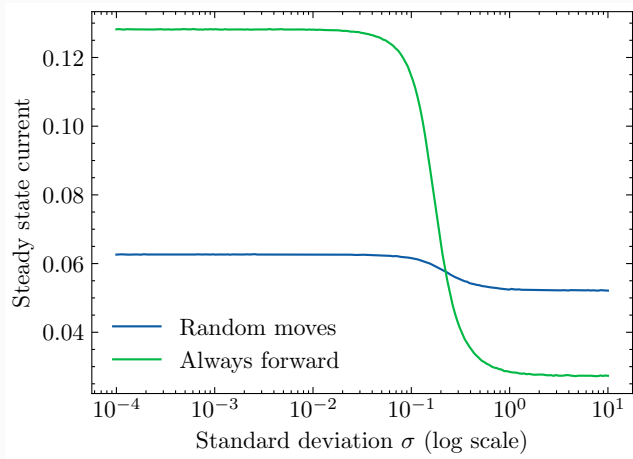Steady state current as a function of $\sigma$ for different system sizes.

Introduction
OOOO

Make it Smart!
OO

Deep Q-Learning
OOOO

Implementation
OOOOO

Baseline
OOO●

Naive Policy
OO

Smarticles
OOOOOOOOOOOOO

Conclusion and Outlook
OO

# Naive Policy

# A Naive Policy

## Just go forward?

- $p_{\text{fwd}} = 1, p_{\text{up/down}} = 0$
- starting in checkerboard configuration
- no messy mixing
- bad for high $\sigma$
- optimum for low $\sigma$?

$$\max \langle J \rangle = \max \left( p_{\text{occ}} \cdot p_{\text{spd}} \cdot p_{\text{fwd}} \cdot p_{\text{emp}} \right)$$
$$\leq \rho \cdot \mu \cdot 1 \cdot \max(p_{\text{emp}})$$
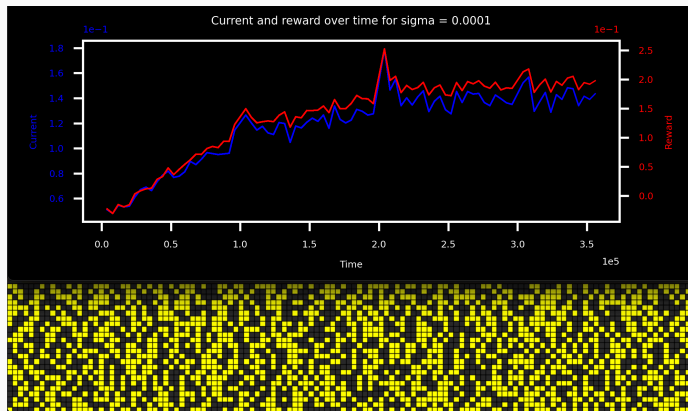
$$\text{and} \quad 0.5 \leq \max(p_{\text{emp}}) < 1$$



Naive policy

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

**Naive Policy**
○●

Smarticles
○○○○○○○○○○○○○

Conclusion and Outlook
○○

# Smarticles

## Simple reward

- Positive reward for going forward
- Negative reward for occupied destination
- Zero reward for going up or down or staying
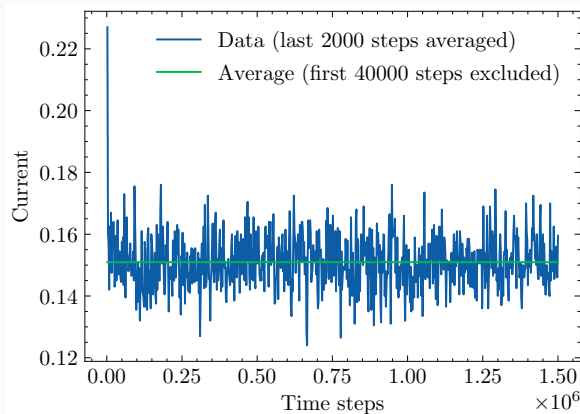- Small negative reward for blocking others



Screenshot of the first smarticle training

## Resulting current

- $\langle J \rangle \approx 0.152$
- $+143\%$ compared to baseline (0.0625)
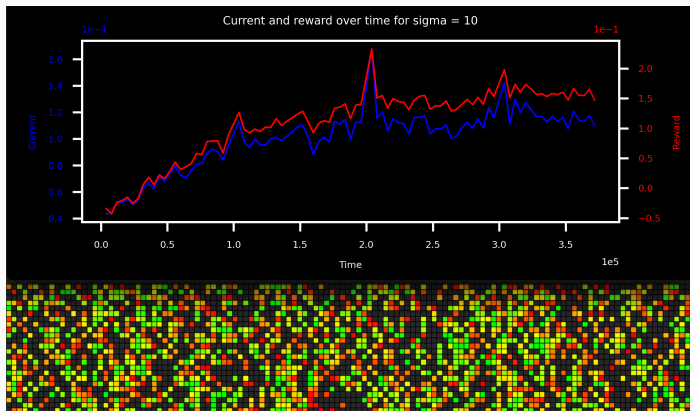- $+22\%$ compared to naive policy (0.125)



Steady state current over time for the trained smarticles.

**Similar reward, $\sigma = 10$**

- Same setup as before
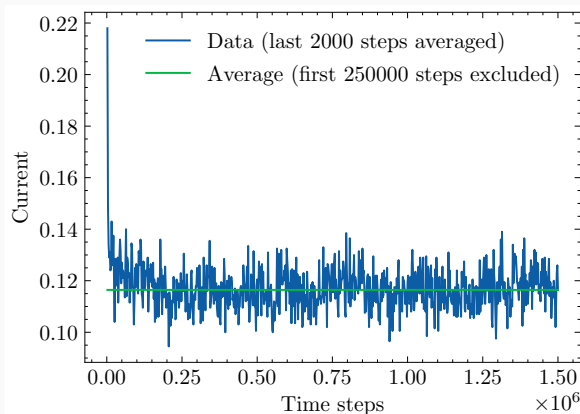- Negative reward for blocking proportional to speed



Screenshot of the second smarticle training

## Resulting current

- $\langle J \rangle \approx 0.117$
- $+125\%$ compared to baseline (0.052)
- Lower than with equal speeds (we will see why)



Steady state current over time for the trained smarticles.

Introduction
OOOO

Make it Smart!
OO

Deep Q-Learning
OOOO

Implementation
OOOOO

Baseline
OOOO

Naive Policy
OO

Smarticles
OOOO●OOOOOOOO

Conclusion and Outlook
OO

What did the smarticles learn? (Playground video)

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

**Smarticles**
○○○○○●○○○○○○

Conclusion and Outlook
○○

## Observation

- So far, the system still looks messy
- Only individual behavior
- Compare highway: Fast lane and slow lane → Could be beneficial also in TASEP

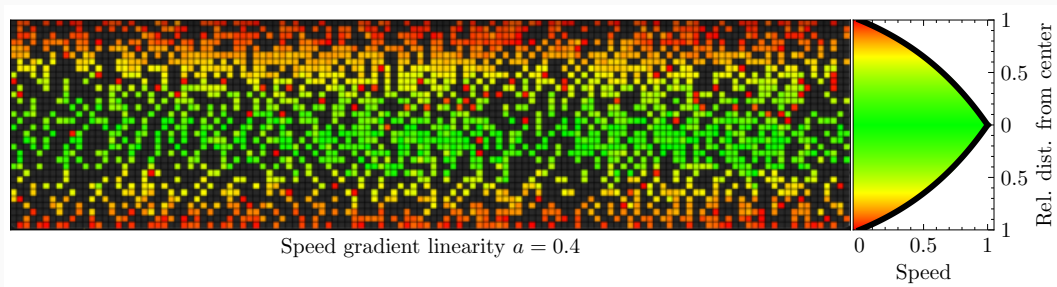messy tasep pic fast lane slow lane picture

## Potential

- Less jamming
- Directed flow without dispersion
- Inhomogeneous density

Introduction
0000

Make it Smart!
OO

Deep Q-Learning
0000

Implementation
00000

Baseline
0000

Naive Policy
OO

Smarticles
000000●000000

Conclusion and Outlook
OO

Approach 1: Hard code the lanes, supply smarticles with lane information

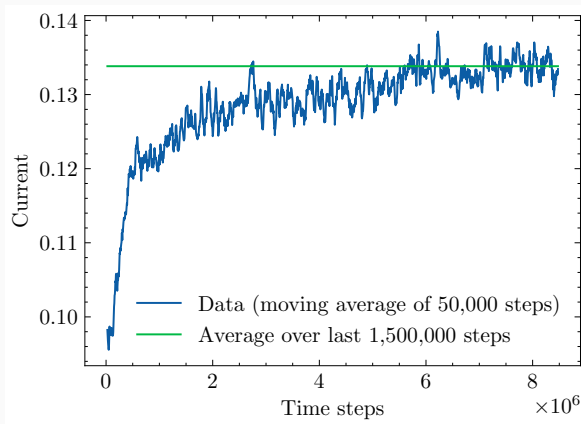$$\Delta y = 1 - a \cdot \left(\frac{1+a}{a}\right)^v + a$$



Speed gradient linearity $a = 0.4$

Learned smarticle behavior with the hard-coded lanes reward function (left), and the speed gradient mapping (right).

Introduction
OOOO

Make it Smart!
OO

Deep Q-Learning
OOOO

Implementation
OOOOO

Baseline
OOOO

Naive Policy
OO

**Smarticles**
OOOOOOO●OOOOO

Conclusion and Outlook
OO

## Resulting current

- $\langle J \rangle \approx 0.134$
- +157% compared to baseline (0.052)
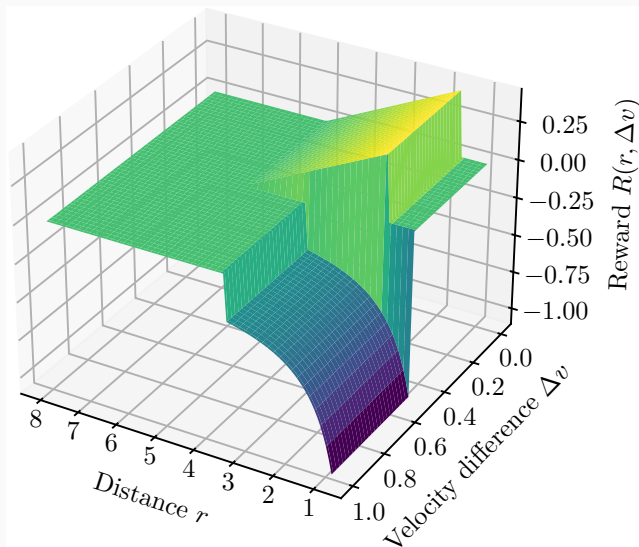- +15% compared to previous policy (0.117)
- long equilibration time



Current over time for the trained smarticles.

Introduction
OOOO

Make it Smart!
OO

Deep Q-Learning
OOOO

Implementation
OOOOO

Baseline
OOOO

Naive Policy
OO

Smarticles
OOOOOOO●OOOOO

Conclusion and Outlook
OO

Approach 2: Let the smarticles learn the lanes themselves from local "interactions"

$$V(r, \Delta v) = \begin{cases} -0.125 \cdot r + 0.625 & \text{if } \Delta v < 0.5 \text{ and } 1.5 < r \leq 5 \\ -0.75 \cdot r^{-1.3} - 0.15 & \text{if } \Delta v \geq 0.5 \text{ and } r \leq 3.5 \\ 0 & \text{otherwise} \end{cases}$$

### Modifications

1. Binary speed distribution
2. Different neural networks

Introduction
0000

Make it Smart!
00

Deep Q-Learning
0000

Implementation
00000

Baseline
0000

Naive Policy
00

Smarticles
000000000●000

Conclusion and Outlook
00

Visual representation of the lane reward function / "potential"

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○●○○

Conclusion and Outlook
○○

Let's see it in action! (lane formation video)

Introduction
0000

Make it Smart!
00

Deep Q-Learning
0000

Implementation
00000

Baseline
0000

Naive Policy
00

**Smarticles**
000000000000●0

Conclusion and Outlook
00

## Resulting current

- $\langle J \rangle \approx 0.112$
- baseline
  $J = 0.4 \cdot (1 - 0.4) \cdot 0.6 \cdot 0.5 = 0.072$
- $+56\%$ compared to baseline
- lower than previous policies
- quicker equilibration
- fluctuations seem higher than they are
- proof of concept, room for improvement



Current over time for the trained smarticles.

Introduction
○○○○

Make it Smart!
○○

Deep Q-Learning
○○○○

Implementation
○○○○○

Baseline
○○○○

Naive Policy
○○

Smarticles
○○○○○○○○○○○○●

Conclusion and Outlook
○○

# Conclusion and Outlook

Introduction
0000

Make it Smart!
00

Deep Q-Learning
0000

Implementation
00000

Baseline
0000

Naïve Policy
00

Smarticles
000000000000

Conclusion and Outlook
○●