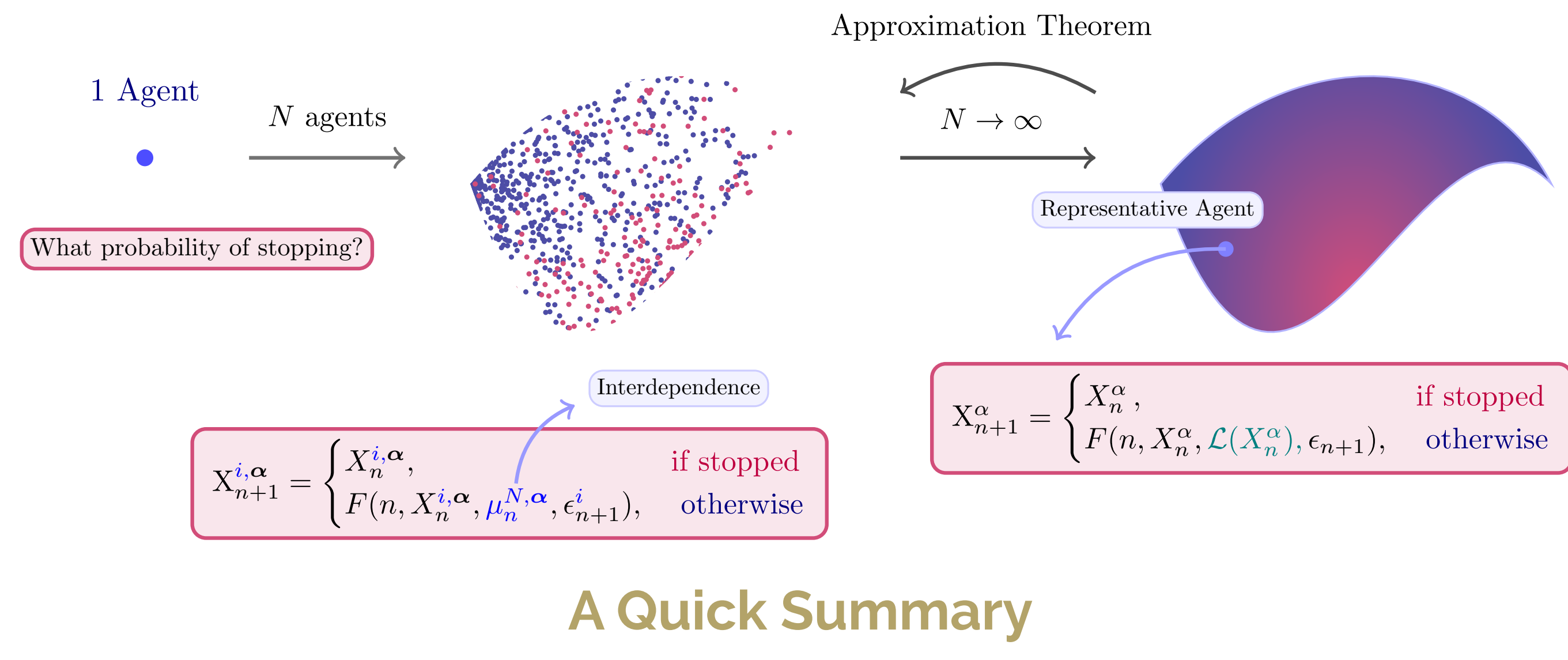


Learning to Stop: Deep Learning for Mean Field Optimal Stopping

Lorenzo Magnino^{*,1} Yuchen Zhu^{*,2} Mathieu Laurière¹
¹NYU Shanghai, ²Georgia Institute of Technology



Target Problem:

- Find practical algorithms for solving Optimal Stopping of Mean-field dynamics of **discrete-time and finite state space**

Importance:

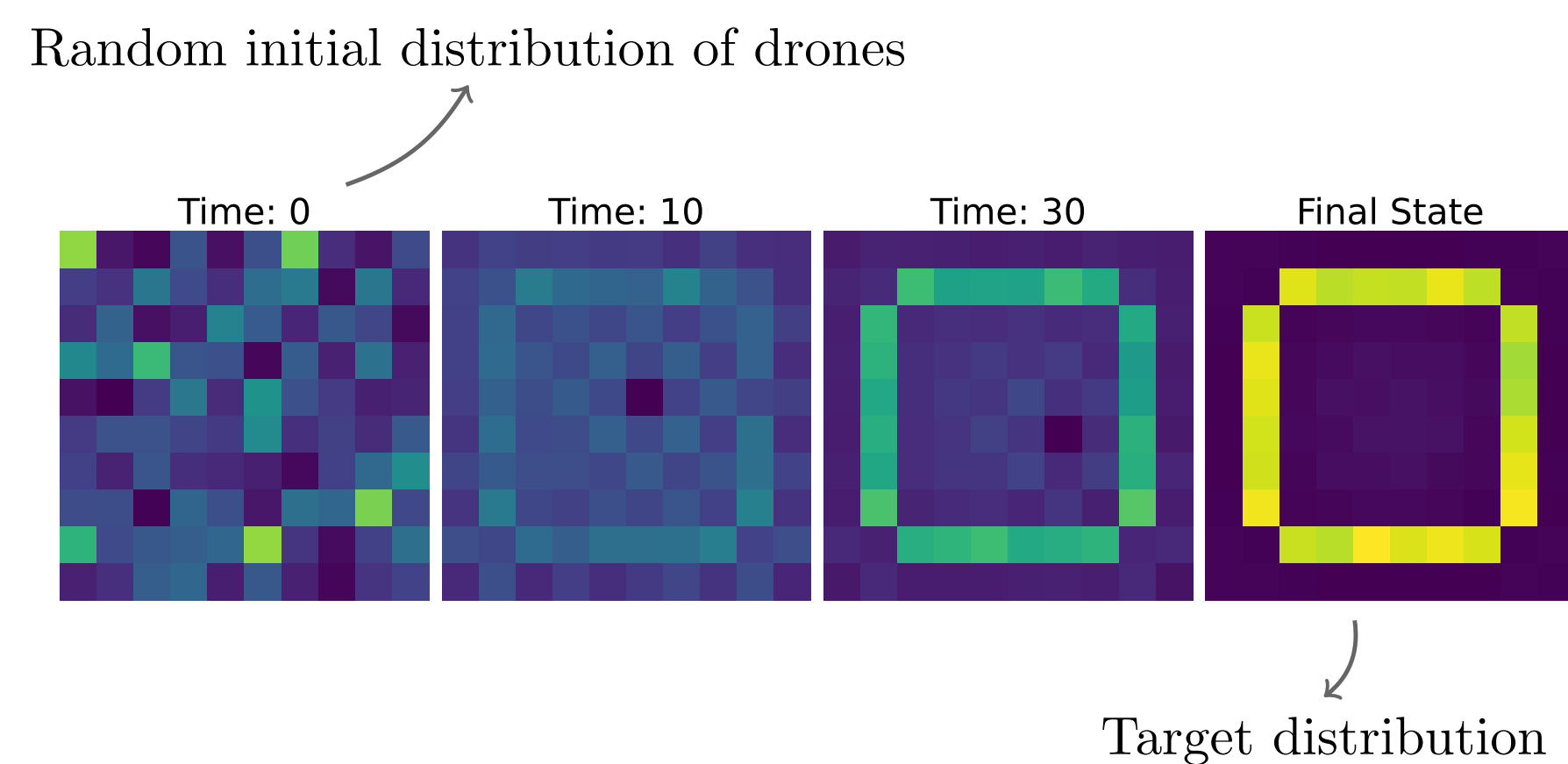
- A proxy for solving Multi-agent OS due to **law of representative agent**
- Applications: option pricing, swarm robotics, etc

Approach:

- Introduce extra **extended state** to signal the decision status of each agent
- Relate Mean-field Optimal Stopping to Mean-field Control to build theory

Results:

- Establish $\mathcal{O}(1/\sqrt{N})$ approximation error between N -agent MAOS and MFOS
- Propose two algos: **Direct Approach** and **Dynamic Programming Principle**
- Consider two type of policies: **synchronous** and **asynchronous**
 - synchronous**: agent decides based on population distribution (suboptimal in many cases)
 - asynchronous** agent decides based on **current state, time** and population



Mean-field Optimal Stopping

We are concerned with finding control α for the cost:

$$\text{Multi-agent: } J^N(\alpha^1, \dots, \alpha^N) = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \Phi(X_{\tau^i}^{i,\alpha}, \mu_{\tau^i}^{N,\alpha}) \right]$$

$$\text{Mean-field: } J(\alpha) = \mathbb{E} \left[\Phi(X_\tau^\alpha, \mathcal{L}(X_\tau^\alpha)) \right].$$

Agent dynamic with **extended state**:

$$\begin{cases} X_0^\alpha \sim \mu_0, & A_0^\alpha = 1 \\ \alpha_n \sim \pi_n(\cdot | X_n^\alpha) = \text{Ber}(p_n(X_n^\alpha)); & A_{n+1}^\alpha = A_n^\alpha(1 - \alpha_n) \\ X_{n+1}^\alpha = \begin{cases} F(n, X_n^\alpha, \mathcal{L}(X_n^\alpha), \epsilon_{n+1}), & \text{if } A_n^\alpha(1 - \alpha_n) = 1 \\ X_n^\alpha, & \text{otherwise.} \end{cases} \end{cases}$$

Mean-field dynamic evolution:

$$\begin{cases} \nu_0^p(x, 0) = 0, & \nu_0^p(x, 1) = \mu_0(x), & \nu_{n+1}^p = \bar{F}(\nu_n^p, p_n), \\ (\bar{F}(\nu, h))(x, a) = \left(\nu(x, 0) + \nu(x, 1)h(x) \right) (1 - a) + \left(\sum_{z \in \mathcal{X}} \nu(z, 1) \left(q_{z,x}^\nu (1 - h(z)) \right) \right) a, \end{cases}$$

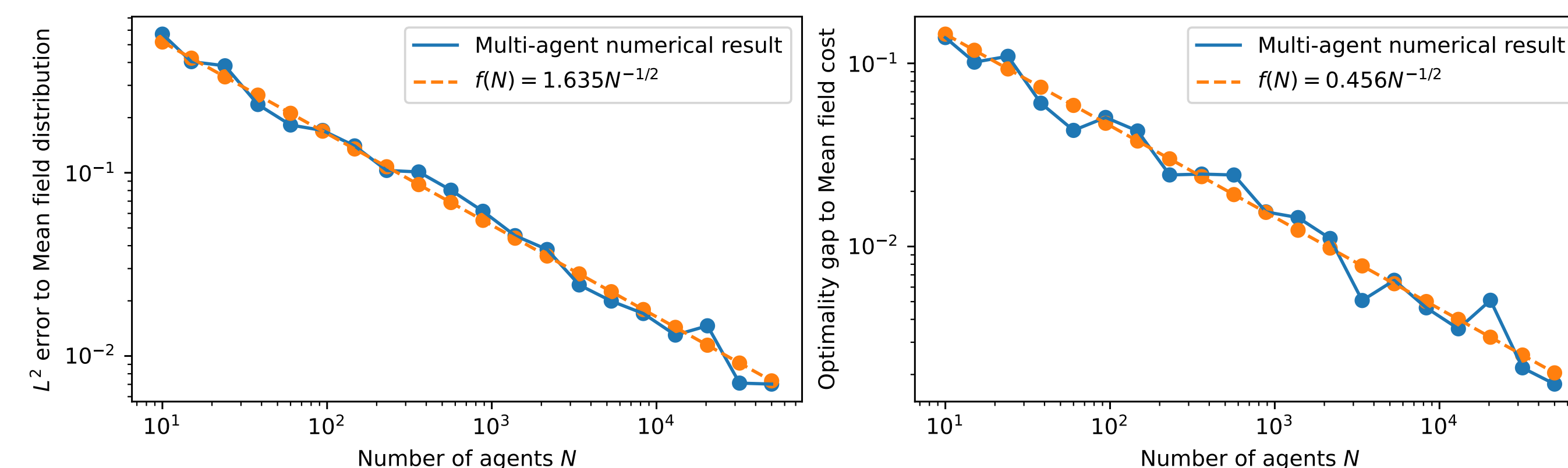
Dynamic Programming Principle:

$$V_n(\nu) := \inf_{p \in \mathcal{P}_{n,T}} J(p(x), \nu) = \inf_{p \in \mathcal{P}_{n,T}} \sum_{m=n}^T \sum_{(x,a) \in \mathcal{S}} \nu_m^{p,\nu,n}(x, a) \Phi(x, \mu_m^{p,\nu,n}) a p_m(x),$$

$$\begin{cases} V_T(\nu) = \sum_{(x,a) \in \mathcal{S}} \nu(x, a) \Phi(x, \nu_X) a, \\ V_n(\nu) = \inf_{h \in \mathcal{H}} \sum_{(x,a) \in \mathcal{S}} \nu(x, a) \Phi(x, \nu_X) a h(x) + V_{n+1}(\bar{F}(\nu, h)), & n < T, \end{cases}$$

ε -Approximation Error: p^* optimal for MFOS, \hat{p} optimal for MAOS (each agent use same policy), $\mathcal{O}(1/\sqrt{N})$ error

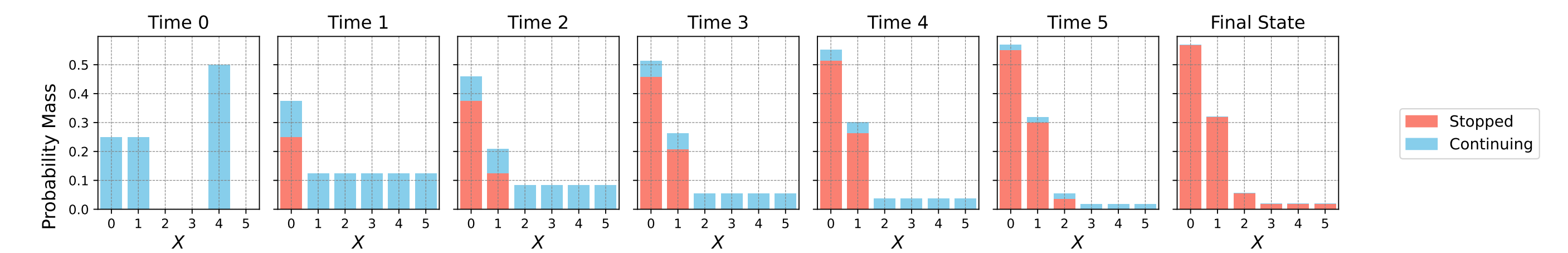
$$J^N(p^*, \dots, p^*) - J^N(\hat{p}, \dots, \hat{p}) \leq 2TL_\Psi(1 + L_p) \left[\frac{|S|}{4\sqrt{N}} \left(\frac{1 - (L_{\bar{F}}(1 + L_p))^T}{1 - (L_{\bar{F}}(1 + L_p))} \right) + (L_{\bar{F}}(1 + L_p))^T \frac{\sqrt{|S| - 1}}{2\sqrt{N}} \right]$$



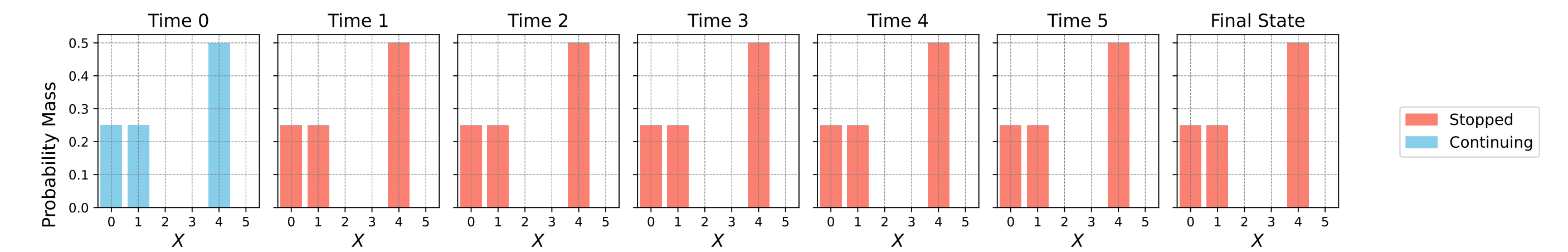
Rolling a Die (for 4 times!)

Setting: Each agent rolls a die and decides whether to stop. If one stops, he pays the cost of the **current number** he had. If one is not satisfied, he can reroll the die (up to 4 time). 25% starts with 1, 25% starts with 2, 50% starts with 5.

Asynchronous: $V^* = 1.6525$, stop when landing on smaller numbers



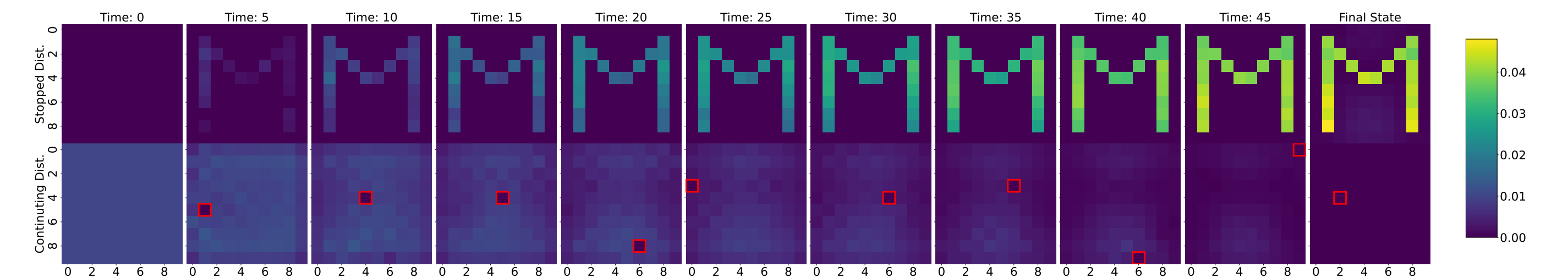
Synchronous: $V^* = 3.25$, stop at beginning



Matching the Letter (on 10 by 10 grid!)

Setting: Each agent is a drone starting randomly somewhere on the grid. For 50 steps, the drone diffuses uniformly to **accessible neighboring positions**. The goal is to finally end up with a distribution of drones that **matches a given letter (like a show!)**. An obstacle shows up in a random position to block the way each time.

Trajectory Snapshots:



Robustness to Initial distributions:

