

# Deep Reinforcement Learning of Marked Temporal Point Processes

Max Planck Institute for Software Systems – Networks and Machine Learning

Utkarsh Upadhyay, Abir De, Manuel Gomez-Rodriguez

{utkarshu,ade,manuelgr}@mpi-sws.org; [Networks-Learning/TPPRL](https://github.com/utkarshu/Networks-Learning/TPPRL)

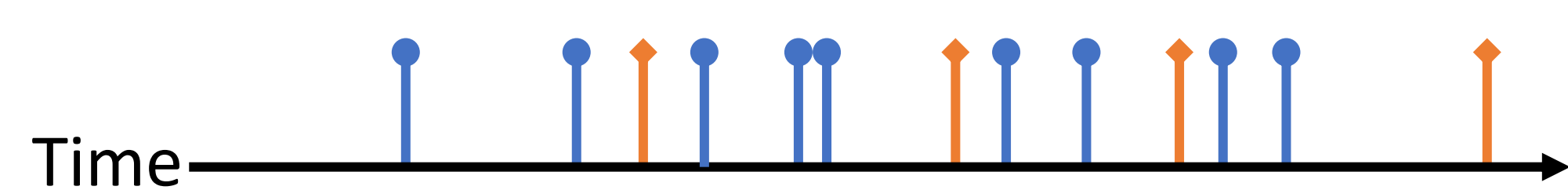


## Introduction

Humans take actions and receive feedback using

- asynchronous and stochastic
  - discrete events in continuous time
- to achieve *their goals*.

### Events as Marked Temporal Point Processes



$t_i \in \mathbb{R}$  : Times  $z_i \in \mathbb{Z}$  : Marks

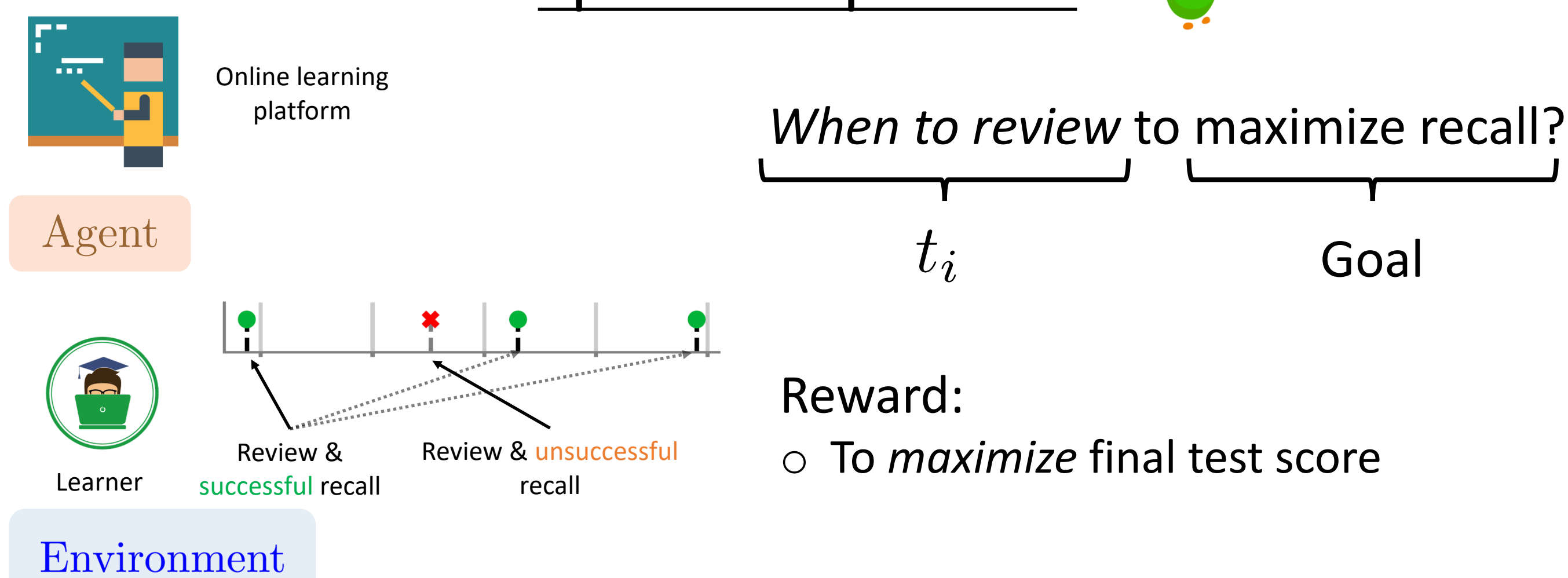
$\mathcal{H}_t = \{e_0 = (t_0, z_0), \dots, e_n = (t_n, z_n)\}$

$\mathcal{H}_t$  : History  $\mathcal{A}_t$  : Actions  $\mathcal{F}_t$  : Feedback

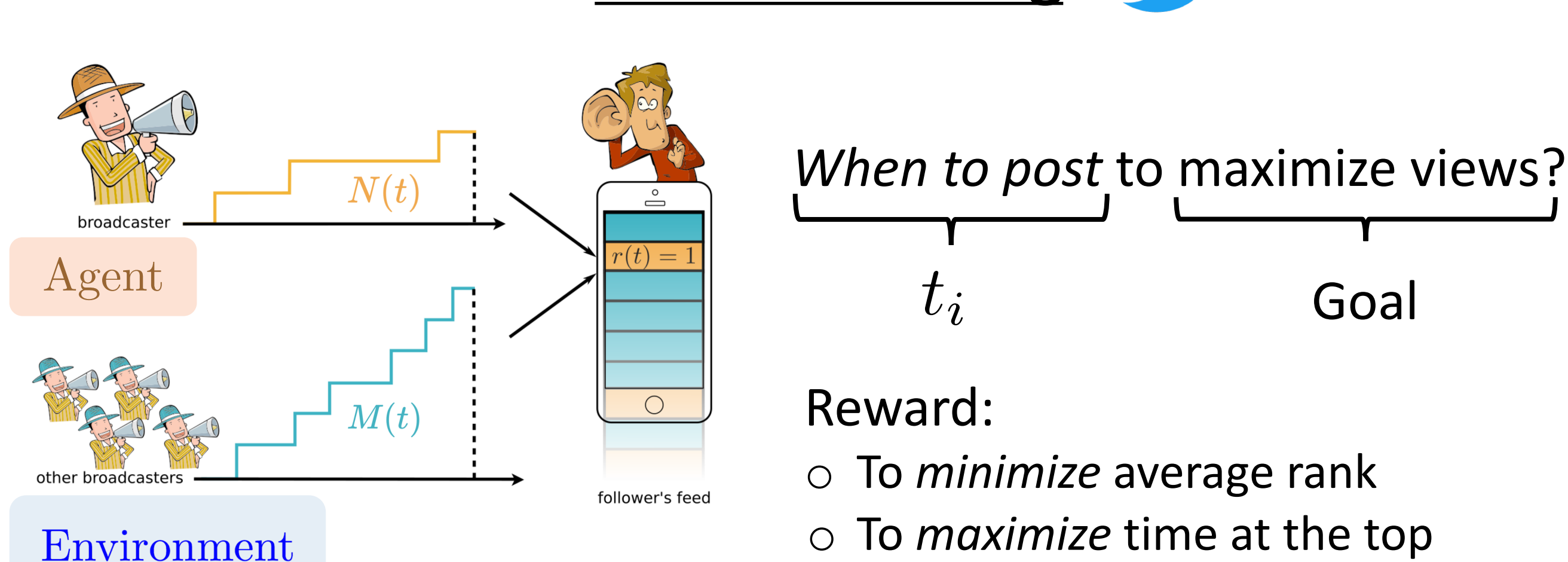
Given the above representation based on MTPPs, can we design online interventions to help them achieve their *goals*?

## Example Applications

### Spaced Repetition

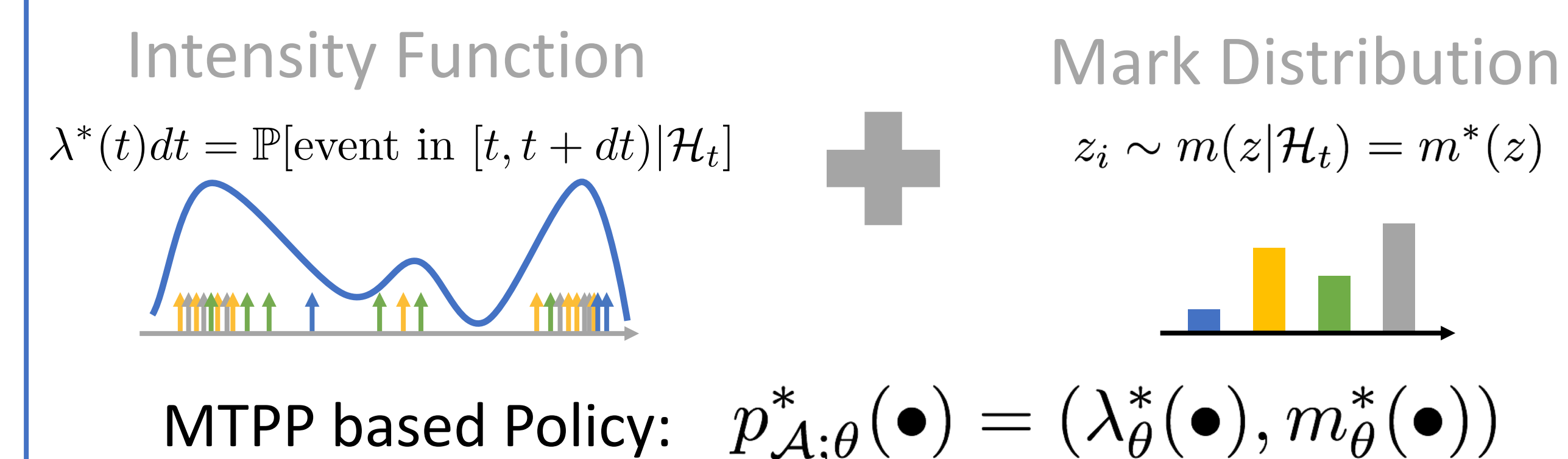


### Viral Marketing

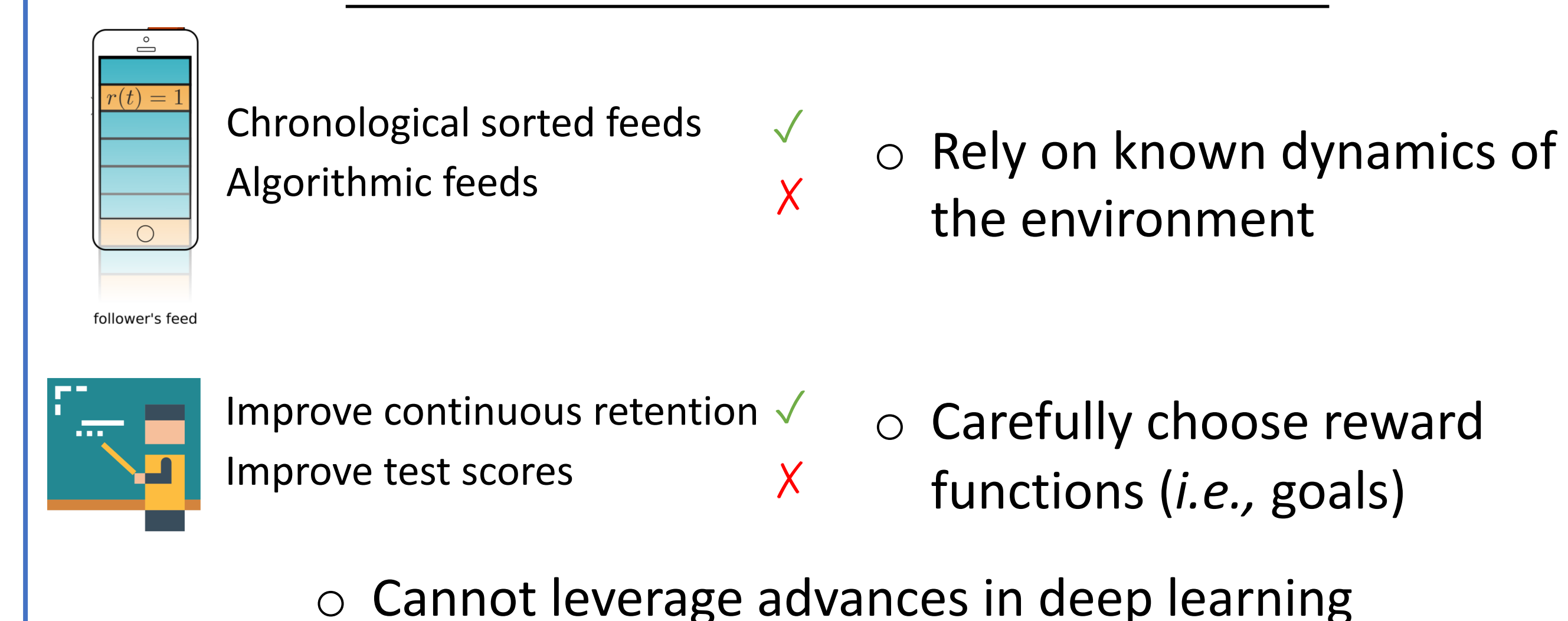


## MTPP based Policy and Control

### Characterization of MTPPs

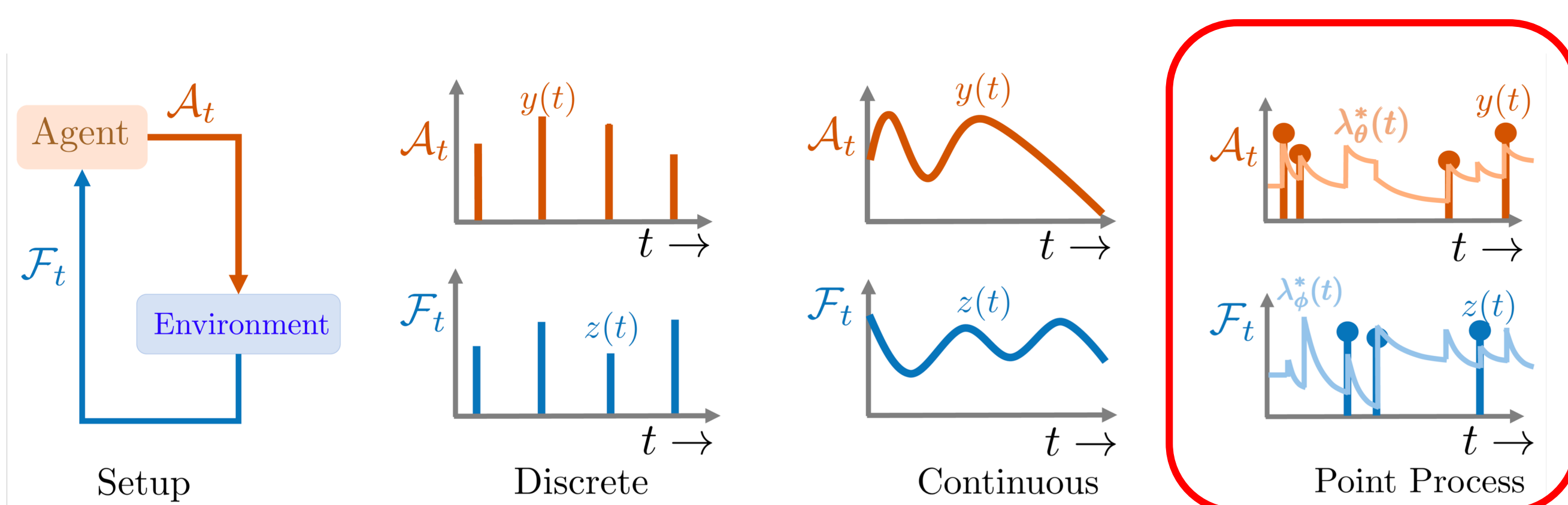


### Previous MTPP Control Methods



## Reinforcement Learning of MTPPs

### Novel RL Setting



### TPPRL: Policy Gradient for MTPPs

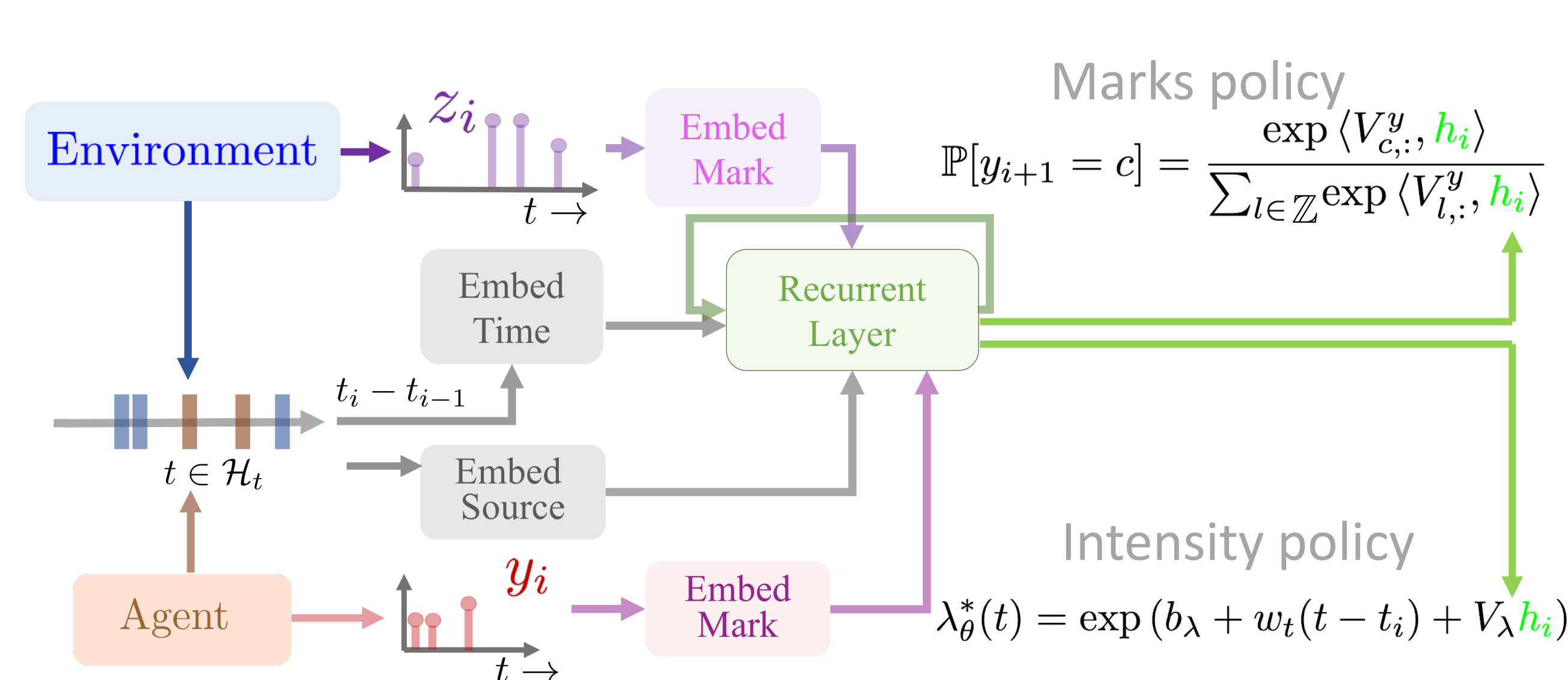
**Aim:** for any reward  $R^*(T)$  maximize  $\mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A},\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F},\phi}^*(\cdot)} [R^*(T)]$

**Solution:** Develop REINFORCE trick for *intensities*.

- Can handle arbitrary reward functions!
- No model for the Feedback process needed!

## Handling Asynchronicity

### Embedding Asynchronous History



### Taking Actions Asynchronously

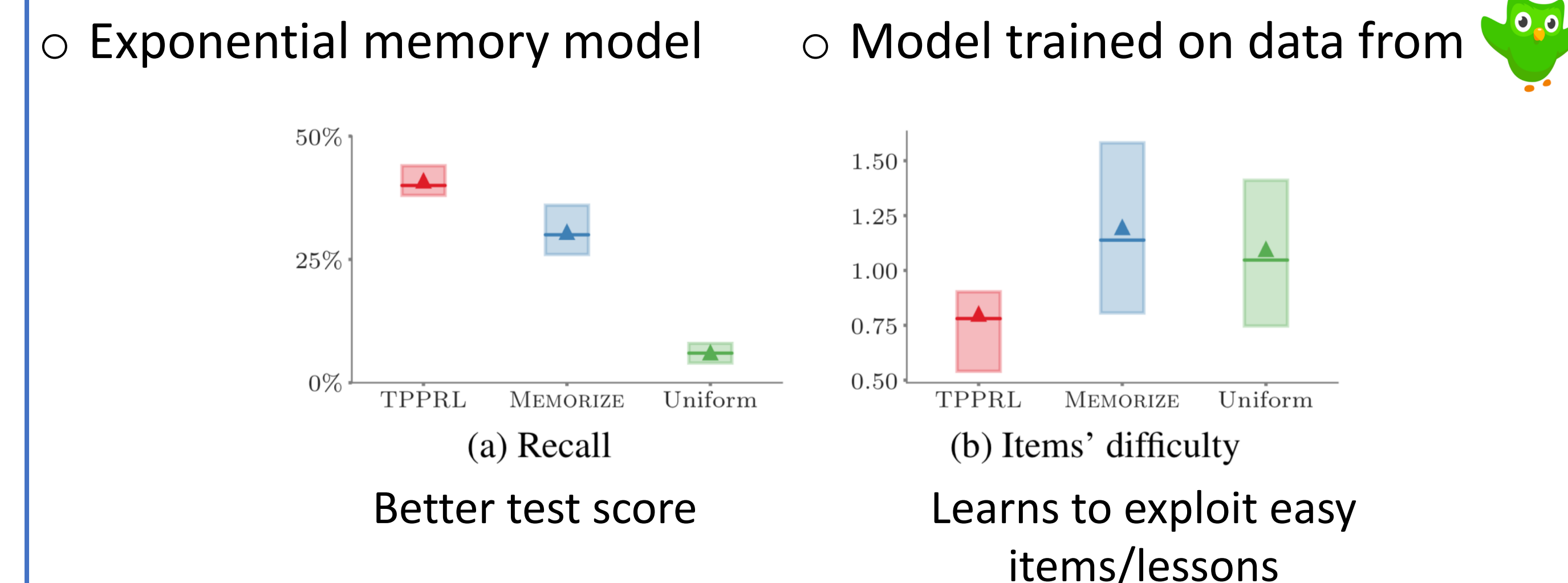
Next event time is sampled from intensity  $\lambda_\theta^*(t)$ , but:

More feedback may arrive in the meanwhile

**Solution:** Efficient & unbiased resampling algorithm.

## Evaluation

### Spaced Repetition



### Viral Marketing

