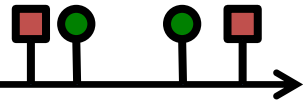


# Reinforcement Learning

## of Marked Temporal Point Processes



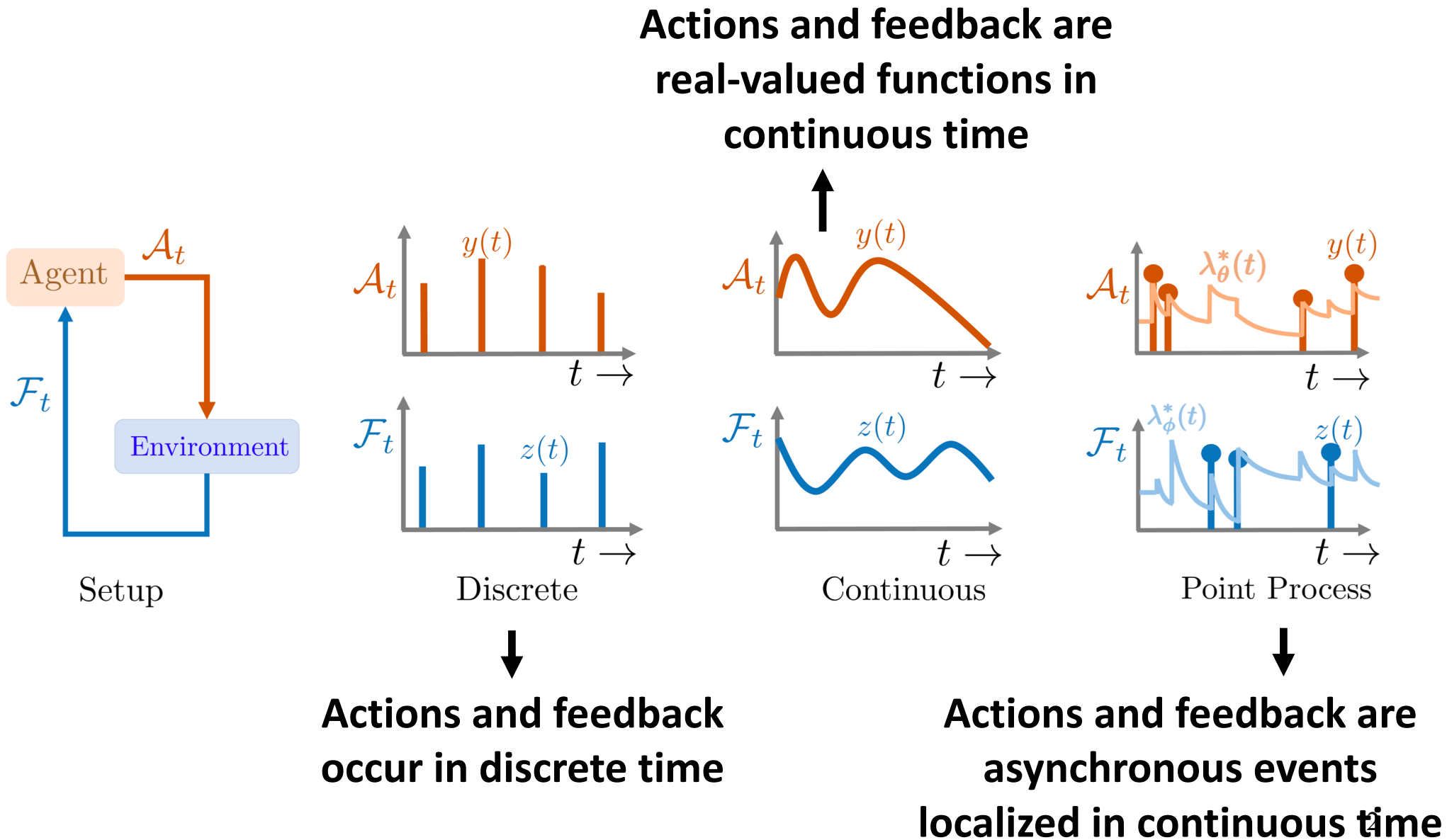
**HUMAN-CENTERED MACHINE LEARNING**

<http://courses.mpi-sws.org/hcml-ws18/>



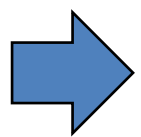
MAX PLANCK INSTITUTE  
FOR SOFTWARE SYSTEMS

# Reinforcement learning on different settings



# Reinforcement learning of marked TPP

If the problem dynamics cannot be expressed using SDEs with jumps or the objective is intractable:



**Reinforcement learning of marked temporal point processes**

→ Policy gradient [Upadhyay, 2018]

→ Policy iteration [Farajtabar et al., 2017]

Similarly as with optimal control:

**Policy is characterized by an intensity function!**

# Reinforcement learning of marked TPP

If the problem dynamics cannot be expressed using SDEs with jumps or the objective is intractable:

→ Next, details on the approach based on policy gradient temporal

→ Policy Iteration [Farajtabar et al., 2017]

Similarly as with optimal control:

**Policy is characterized by an intensity function!**

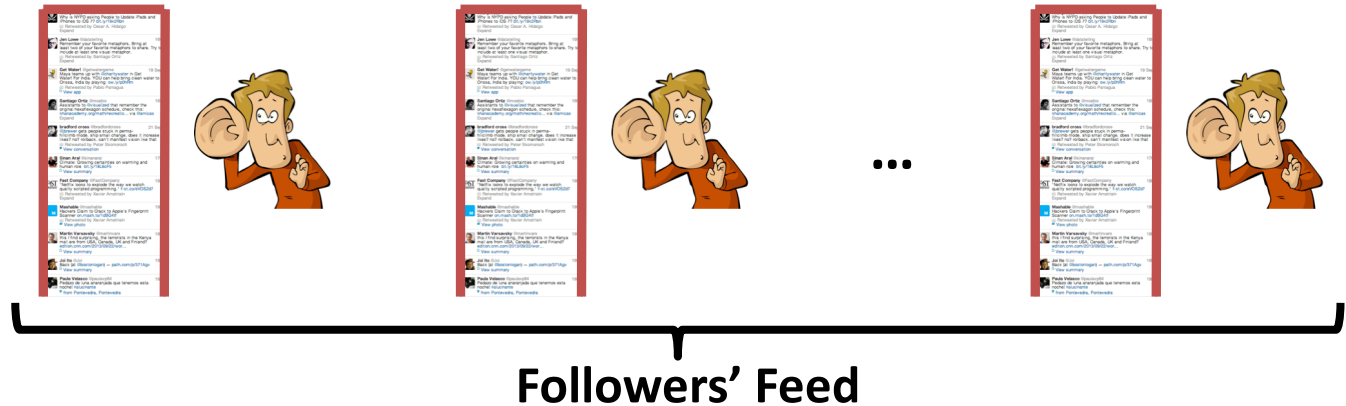
# Viral marketing

## Agent



Social media user

## Environment



## Forbes

For Brands And PR: When Is The Best Time To Post On Social Media?

THE HUFFINGTON POST

The Best Times to Post on Social Media

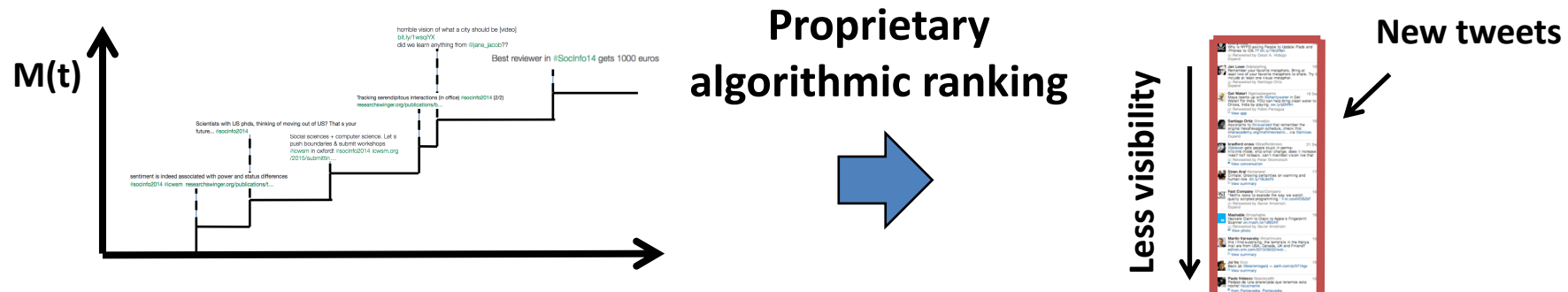
**When to post to maximize views or likes?**

$$\mu_i(t) = u(t) \rightarrow N_i(t)$$

Design (optimal)  
posting intensity

Marks (feedback) given  
by environment

# Visibility dynamics are unknown



However, one may have access to quality metrics

Manuel @autreche  
Three days before the #nips2018 deadline, there already 888 submissions! :)

Reach a bigger audience  
Get more engagements by promoting this Tweet!

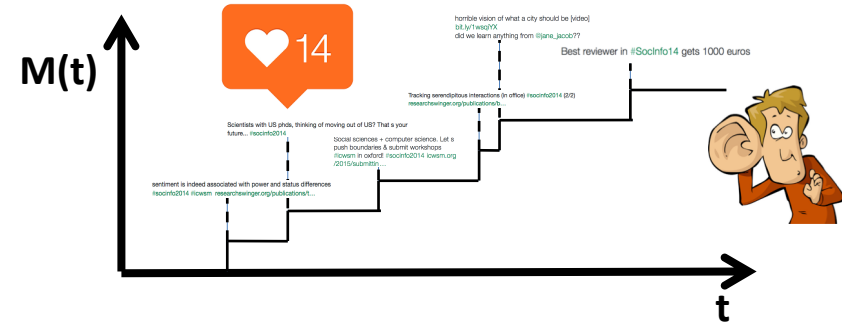
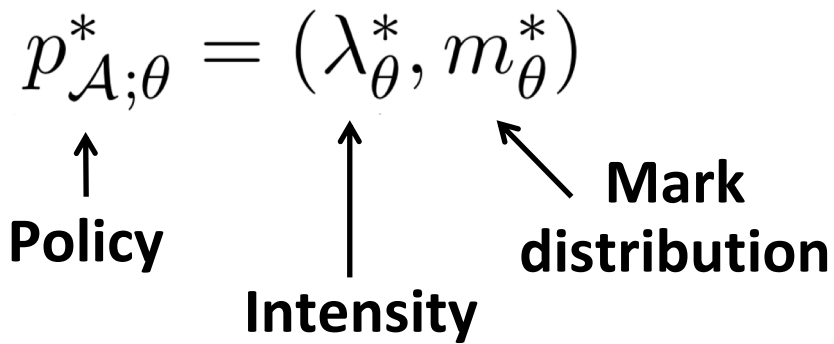
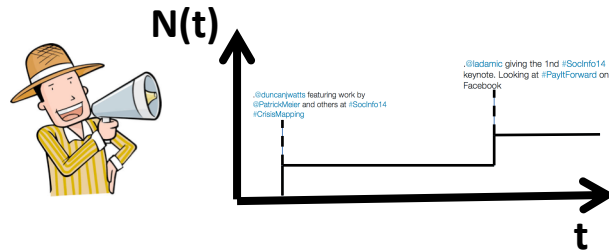
Get started

Impressions	1,096
Total engagements	15
Detail expands	5
Profile clicks	4
Likes	3
Hashtag clicks	3

Key idea:

Think of these metrics as rewards in a reinforcement learning setting!

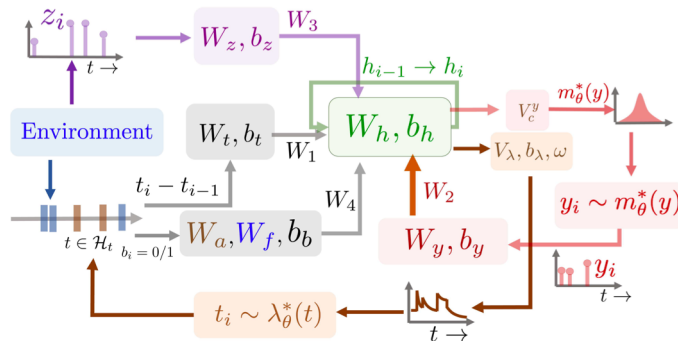
# Broadcasters and feedback



$$p_{\mathcal{F};\phi}^* = (\lambda_{\phi}^*, m_{\phi}^*)$$

We do not know the *feedback* distribution but we can *sample* from it...

## Parametrized using RNNs



Manuel @autreche  
Three days before the #nips2018 deadline, there already 888 submissions! ;]

Reach a bigger audience  
Get more engagements by promoting this Tweet!

[Get started](#)

Impressions	1,096
Total engagements	15
Detail expands	5
Profile clicks	4
Likes	3
Hashtag clicks	3

...and measure **quality metrics** (rewards)

# What is the goal in reinforcement learning?

We aim to maximize the average reward in a time window  $[0, T]$ :

$$\text{maximize}_{p_{\mathcal{A};\theta}^*(\cdot)} \underbrace{\mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T)]}_{\substack{\text{Actions and} \\ \text{environment are} \\ \text{asynchronous!}}} \quad \uparrow \quad \substack{\text{Reward} \\ \text{(Cumulative)}}$$

$J(\theta)$

Connection to optimal control:

$$J(r(t), \lambda(t), t) = \min_{u(t, t_f)} \mathbb{E}_{(N, M)(t, t_f)} \left[ \phi(r(t_f)) + \int_t^{t_f} \ell(r(\tau), u(\tau)) d\tau \right]$$



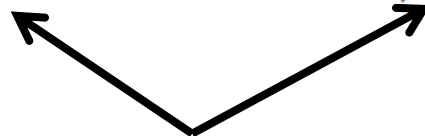
# Policy gradient

We use gradient descent to improve the policy, i.e., the intensity, over time:

$$\theta_{l+1} = \theta_l + \alpha_l \nabla_{\theta} J(\theta) |_{\theta=\theta_l}$$

We need to compute the gradient of an average. But the average depends on the parameters!

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T)]$$

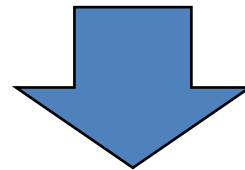


Parameters!

# Reinforce trick to compute gradient

The reinforce trick allows us to overcome this implicit dependence:

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T)]$$



Appendix A in  
Upadhyay et al., 2018

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T) \nabla_{\theta} \log \mathbb{P}_{\theta}(\mathcal{A}_T)]$$

Parameters!

# Likelihood of action events

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T) \nabla_{\theta} \log \mathbb{P}_{\theta}(\mathcal{A}_T)]$$

Likelihood of posts by our broadcaster!

$$\mathbb{P}(\mathcal{A}_T) := \left( \prod_{e_i \in \mathcal{A}_T} \lambda_{\theta}^*(t_i) \right) \exp \left( - \int_0^T \lambda_{\theta}^*(s) ds \right)$$

The key remaining question is how to  
parametrize the intensity  $\lambda_{\theta}^*(t)$



Parameters & functional form!

# Policy parametrization

Parameters

Output layer:

$$\lambda_{\theta}^*(t) = \exp \left( b_{\lambda} + w_t(t - t') + V_{\lambda} h_i \right)$$

Last time the broadcaster posted  
↓

Hidden layer:

$$h_i = \tanh \left( W_h h_{i-1} + W_1 \tau_i + W_4 b_i + b_h \right)$$

Input layer:

$$\tau_i = W_t(t_i - t_{i-1}) + b_t$$

↓            ↓  
i-th event   (i-1)-th event

$$b_i = W_a(1 - e_i) + W_f e_i + b_b$$

$e_i = 0$  action  
 $e_i = 1$  feedback

[Upadhyay et al., 2018]

# Sampling from the policy

$$\lambda_{\theta}^*(t) = \exp(b_{\lambda} + w_t(t - t') + \mathbf{V}_{\lambda} \mathbf{h}_i)$$

**The intensity can increase or decrease every time an event by the other broadcasters take place:**

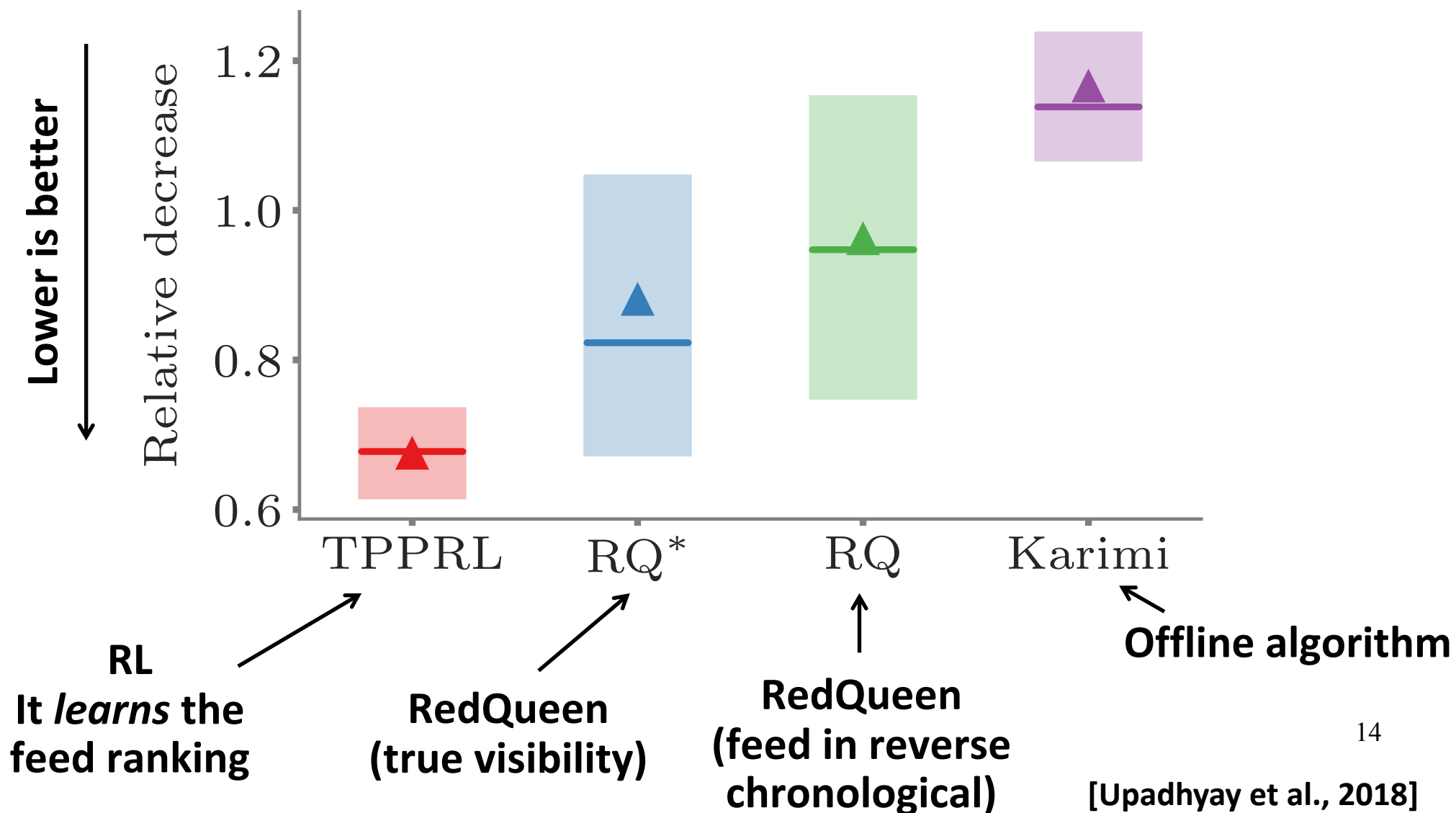
- We cannot apply just superposition**
- We can use inversion sampling: The CDF is a function by parts, where each part is defined once an event by the other broadcasters happens**

Appendix C in  
Upadhyay et al., 2018

[Upadhyay et al., 2018]

# Average Rank

$$R^*(T) = \int_0^T r(t) dt$$



# Time at the top

$$R^*(T) = \int_0^T \mathbb{I}(r(t) < 1) dt$$

