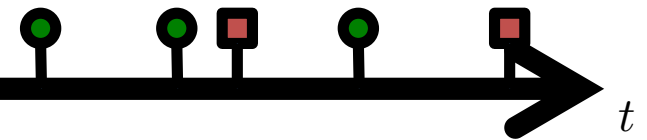


Learning with Temporal Point Processes



RL & Control

Manuel Gomez Rodriguez

Max Planck Institute for Software Systems

Outline of the Seminar

TEMPORAL POINT PROCESSES (TPPs): INTRO

1. Intensity function
2. Basic building blocks
3. Superposition
4. Marks and SDEs with jumps

MODELS & INFERENCE

1. Modeling event sequences
2. Clustering event sequences
3. Capturing complex dynamics
4. Causal reasoning on event sequences

RL & CONTROL

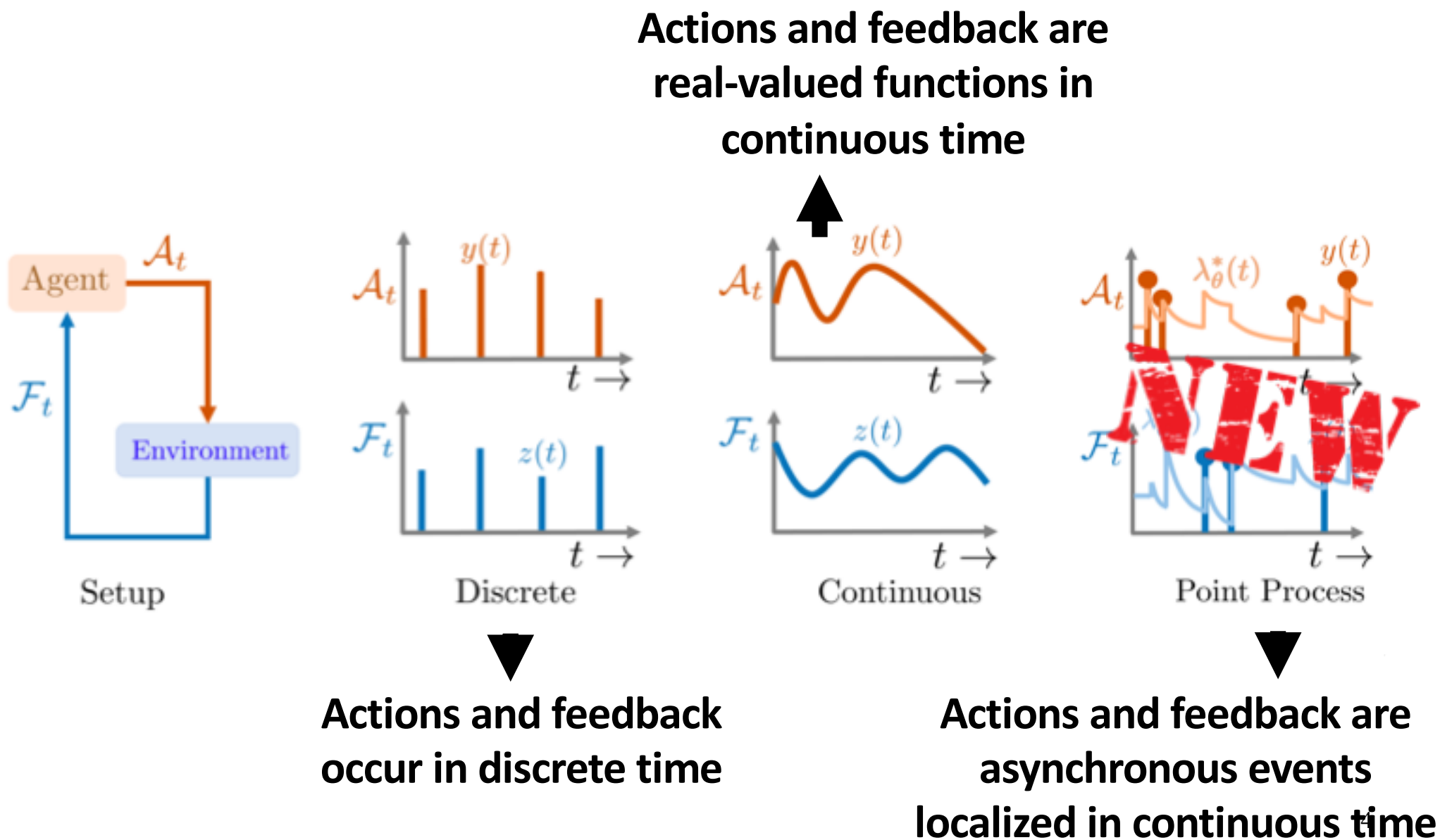
1. Marked TPPs: a new setting
2. Stochastic optimal control
3. Reinforcement learning

Next

RL and Control

- 1. Marked TPP: a new setting**
2. Stochastic optimal control
3. Reinforcement learning

MTPP: a new setting for control & RL



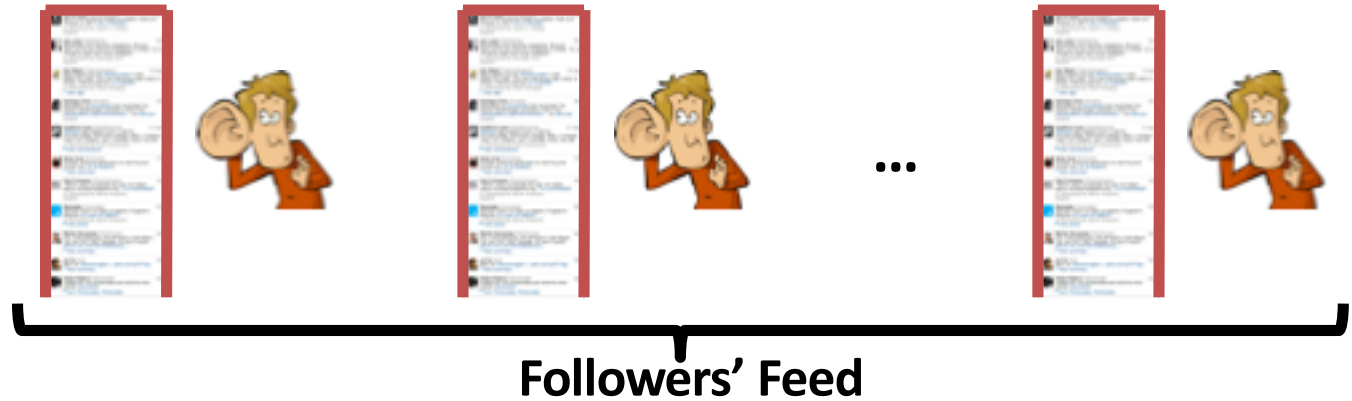
Example I: Viral marketing

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?

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

Design (optimal)
posting intensity

Marks (feedback) given
by environment

Example II: Spaced repetition

Agent

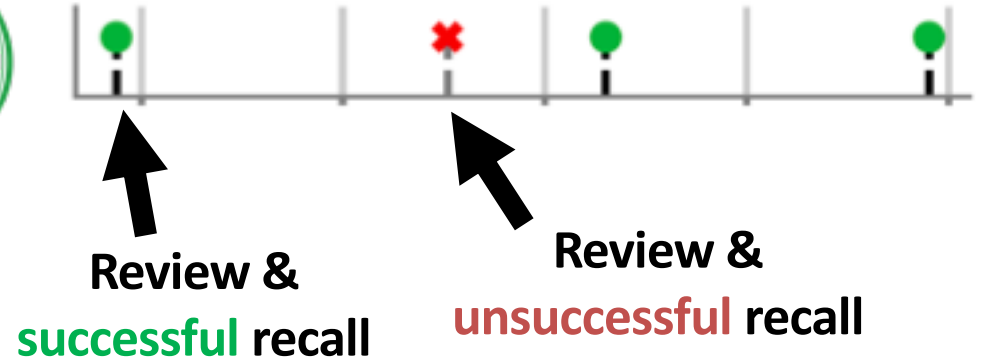


Online learning platform

Environment



Learner



When to review to maximize recall probability?

$$\lambda_i(t) \rightarrow N_i(t)$$

Design (optimal)
reviewing intensities

Marks

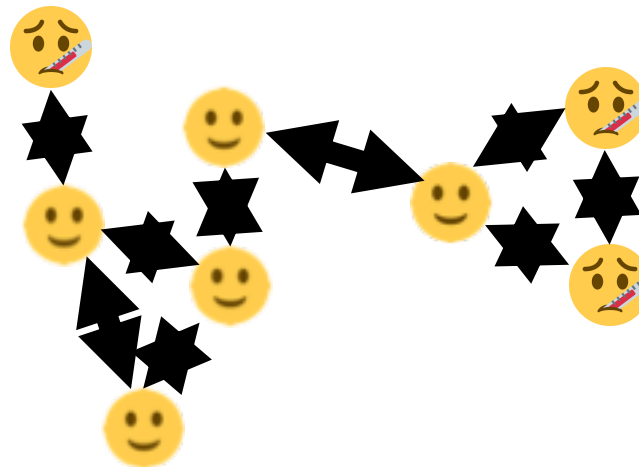
Example III: Suppressing epidemics

Agent



Health policy
(Resource allocation)

Environment



Population (social network)

Who to treat and when to reduce infections?

$$\lambda_i(t) \blacktriangleright N_i(t)$$

Design (optimal)
treatment intensities

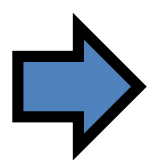
Marks

RL & Control

1. Marked TPP: a new setting for control
- 2. Stochastic optimal control**
3. Reinforcement learning

Stochastic optimal control of SDEs with jumps

If the problem dynamics can be expressed using SDEs with jumps:



Optimal control of marked temporal
po



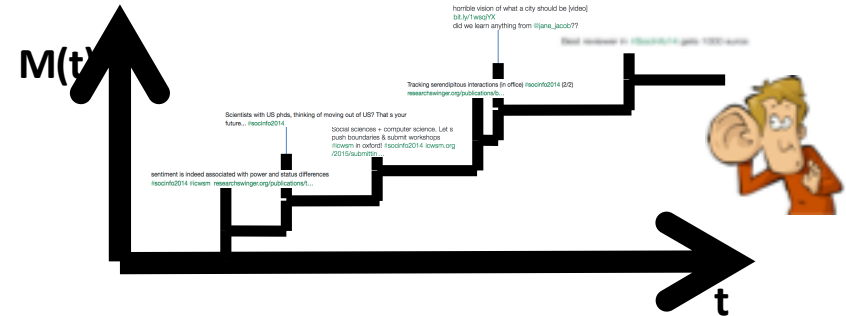
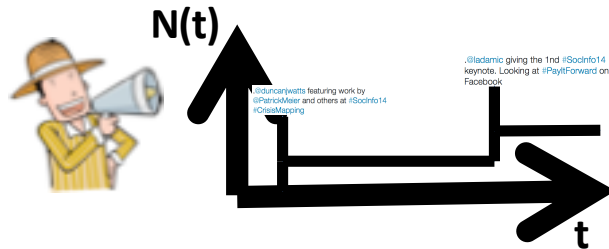
Next, details on one
approach to the when to
post problem

Kim et al. 2018;

Key idea:

Policy is characterized by an intensity
function!

Broadcasters and feeds



$$\mathbb{E}[dN(t)|\mathcal{H}(t)] = \underbrace{\mu(t)} dt \quad \Rightarrow$$

$$\mathbb{E}[dM(t)|\mathcal{H}(t)] = \underbrace{\gamma(t)} dt$$

Policy \Rightarrow

Broadcaster intensity function (tweets / hour)

$$A^T \mu(t)$$

Feed intensity function (tweets / hour)

Given a broadcaster i and her followers \Rightarrow

$$M_{\setminus i}(t) = A^T N(t) - A_i N_i(t)$$

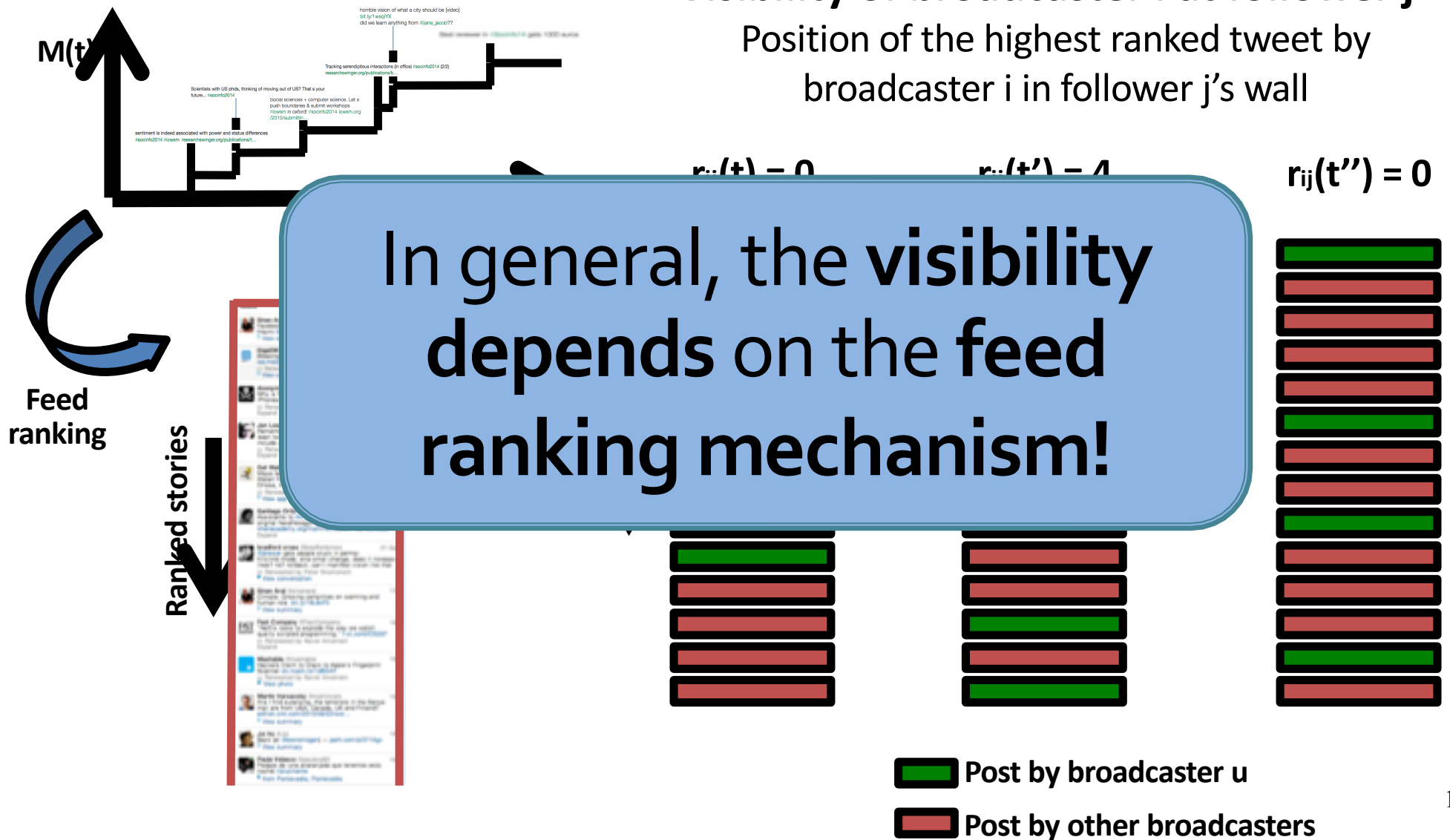
$$\gamma_{j \setminus i}(t) = \gamma_j(t) - \mu_i(t)$$

Feed due to other broadcasters

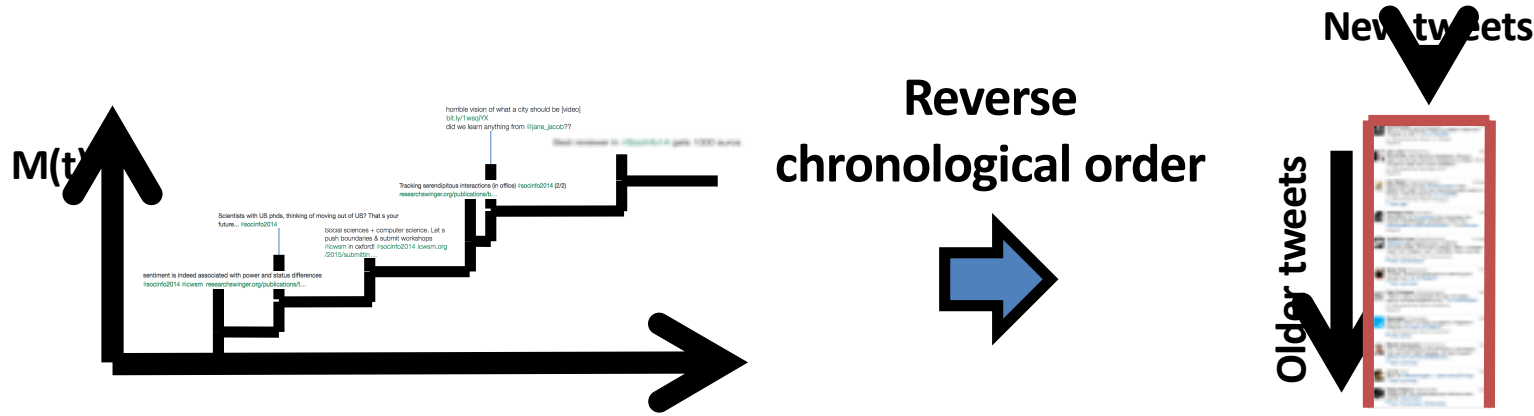
Definition of visibility function

Visibility of broadcaster i at follower j

Position of the highest ranked tweet by broadcaster i in follower j 's wall



Visibility dynamics in a FIFO feed (I)



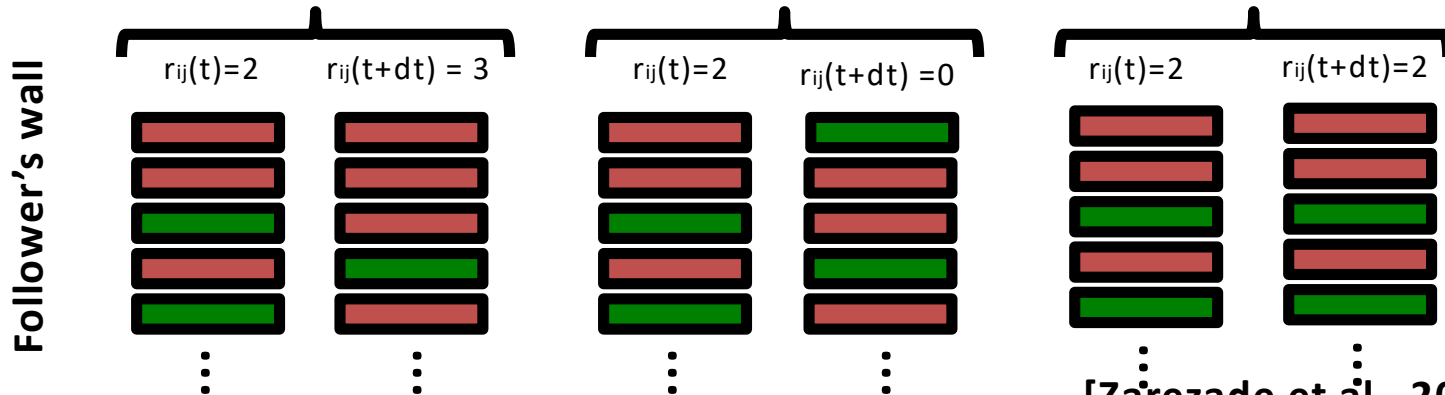
$$r_{ij}(t + dt) = \underbrace{(r_{ij}(t) + 1)}_{\text{Rank at } t+dt} dM_{j \setminus i}(t) \underbrace{(1 - dN_i(t))}_{\text{Other broadcasters post a story and broadcaster } i \text{ does not post}} + 0 + r_{ij}(t) \underbrace{(1 - dM_{j \setminus i}(t))}_{\text{Broadcaster } i \text{ posts a story and other broadcasters do not post}} \underbrace{(1 - dN_i(t))}_{\text{Nobody posts a story}}$$

Rank at $t+dt$

Other broadcasters post a story and broadcaster i does not post

Broadcaster i posts a story and other broadcasters do not post

Nobody posts a story



Visibility dynamics in a FIFO feed (II)

$$r_{ij}(t + dt) = (r_{ij}(t) + 1)dM_{j \setminus i}(t)(1 - dN_i(t)) + 0 + r_{ij}(t)(1 - dM_{j \setminus i}(t))(1 - dN_i(t))$$



Zero-one law $dN_i(t)dM_{j \setminus i}(t) = 0$

$$dr_{ij}(t) = -r_{ij}(t) dN_i(t) + dM_{j \setminus i}(t)$$



$$r_{ij}(t + dt) - r_{ij}(t)$$

Broadcaster *i*
posts a story

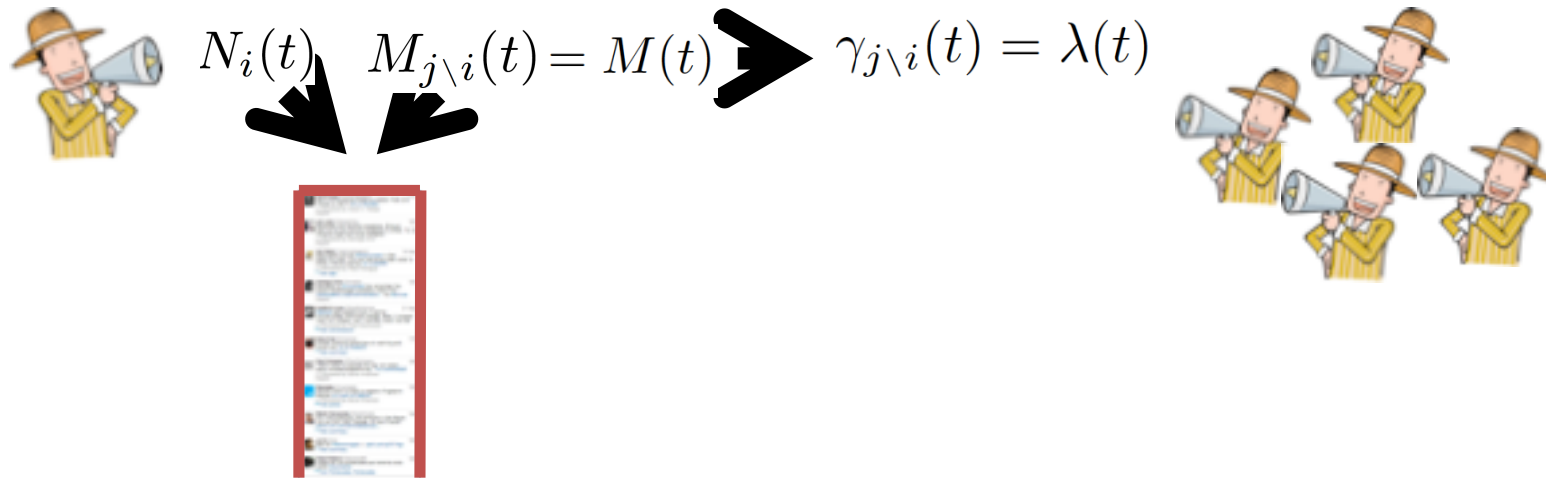
Other broadcasters
posts a story

Stochastic
differential equation
(SDE) with jumps

OUR GOAL:

Optimize $r_{ij}(t)$ over time, so that it is small, by controlling $dN_i(t)$ through the intensity $\mu_i(t)$

Feed dynamics



We consider a **general intensity:**

(e.g. Hawkes, inhomogeneous Poisson)

$$\lambda^*(t) = \underbrace{\lambda_0(t)}_{\text{Deterministic arbitrary intensity}} + \underbrace{\alpha \int_0^t g(t-s) dN(s)}_{\text{Stochastic self-excitation}}$$

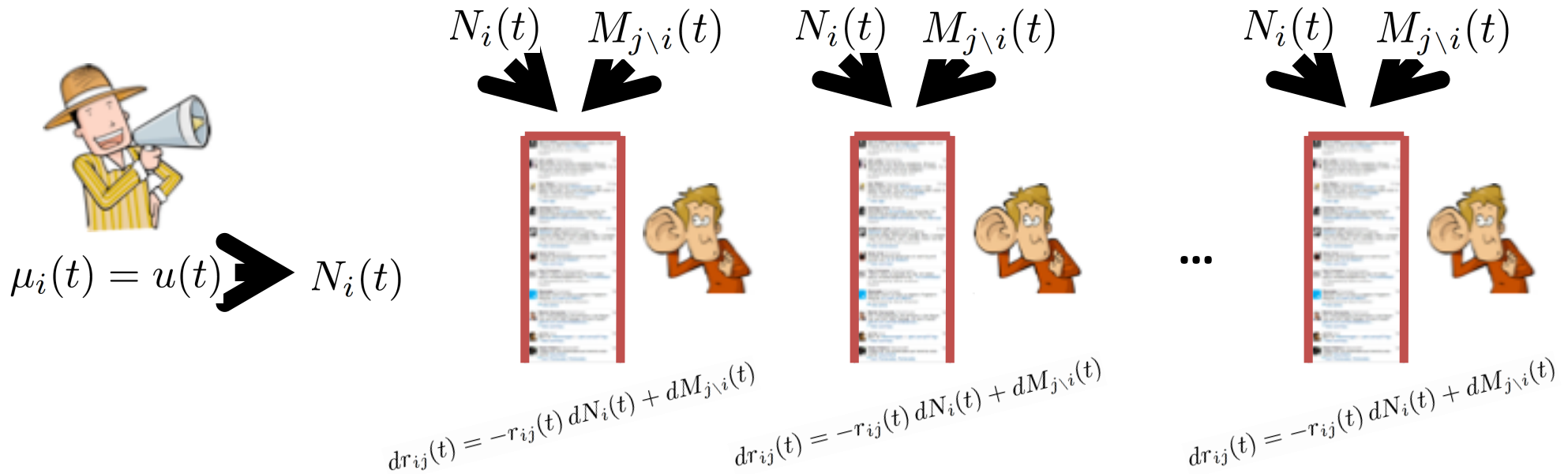


Jump stochastic differential equation (SDE)

$$\left\{ \begin{aligned} d\lambda^*(t) &= [\lambda_0'(t) + w\lambda_0(t) - w\lambda^*(t)] dt + \alpha dN_i(t) \end{aligned} \right.$$

[Zaregade et al., 2017 & 2018]

The when-to-post problem



Optimization problem

minimize $u(t_0, t_f)$ $\mathbb{E}_{(N_i, M_{\setminus i})(t_0, t_f)}$

subject to $u(t) \geq 0 \quad \forall t \in (t_0, t_f]$

Dynamics defined by Jump SDEs

$dr(t) = -r(t) dN(t) + dM(t)$

$d\lambda(t) = [\lambda'_0(t) + w\lambda_0(t) - w\lambda(t)] dt + \alpha dM(t)$

Terminal penalty

Nondecreasing loss

$\left[\phi(\mathbf{r}(t_f)) + \int_{t_0}^{t_f} \ell(\mathbf{r}(\tau), u(\tau)) d\tau \right]$

[Zareezade et al., 2017 & 2018]

Bellman's Principle of Optimality

Lemma. The optimal cost-to-go satisfies Bellman's Principle of Optimality

$$J(r(t), \lambda(t), t) = \min_{u(t, t+dt)} \mathbb{E} [J(r(t+dt), \lambda(t+dt), t+dt)] + \ell(r(t), u(t)) dt$$



$$J(r(t+dt), \lambda(t+dt), t+dt) = J(r(t), \lambda(t), t) + dJ(r(t), \lambda(t), t)$$

$$0 = \min_{u(t, t+dt)} \mathbb{E} [dJ(r(t), \lambda(t), t)] + \ell(r(t), u(t)) dt$$



$$\begin{aligned} dr(t) &= -r(t) dN(t) + dM(t) \\ d\lambda(t) &= [\lambda'_0(t) + w\lambda_0(t) - w\lambda(t)] dt + \alpha dM(t) \end{aligned}$$

Hamilton-Jacobi-Bellman (HJB) equation



Partial differential equation in J (with respect to r , λ and t)¹⁶

[Zaregade et al., 2017 & 2018]

Solving the HJB equation

Consider a quadratic loss

$$\ell(r(t), u(t)) = \frac{1}{2} s(t) r^2(t) + \frac{1}{2} q u^2(t)$$

Favors some periods of times
(e.g., times in which the follower is
online)

Trade-offs visibility and number
of broadcasted posts

Then, it can be shown that the optimal intensity is:

$$\begin{aligned} u^*(t) &= q^{-1} [J(r(t), \lambda(t), t) - J(0, \lambda(t), t)] \\ &= \sqrt{s(t)/q} r(t) \end{aligned}$$

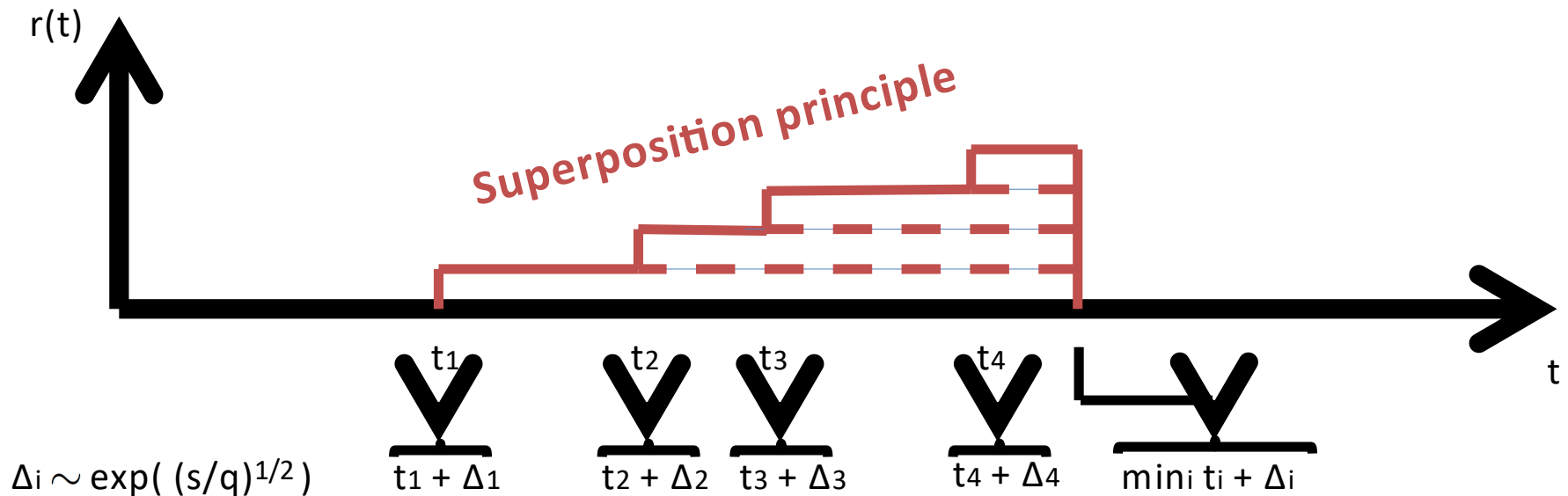
**It only depends on the
current visibility!**



The RedQueen algorithm

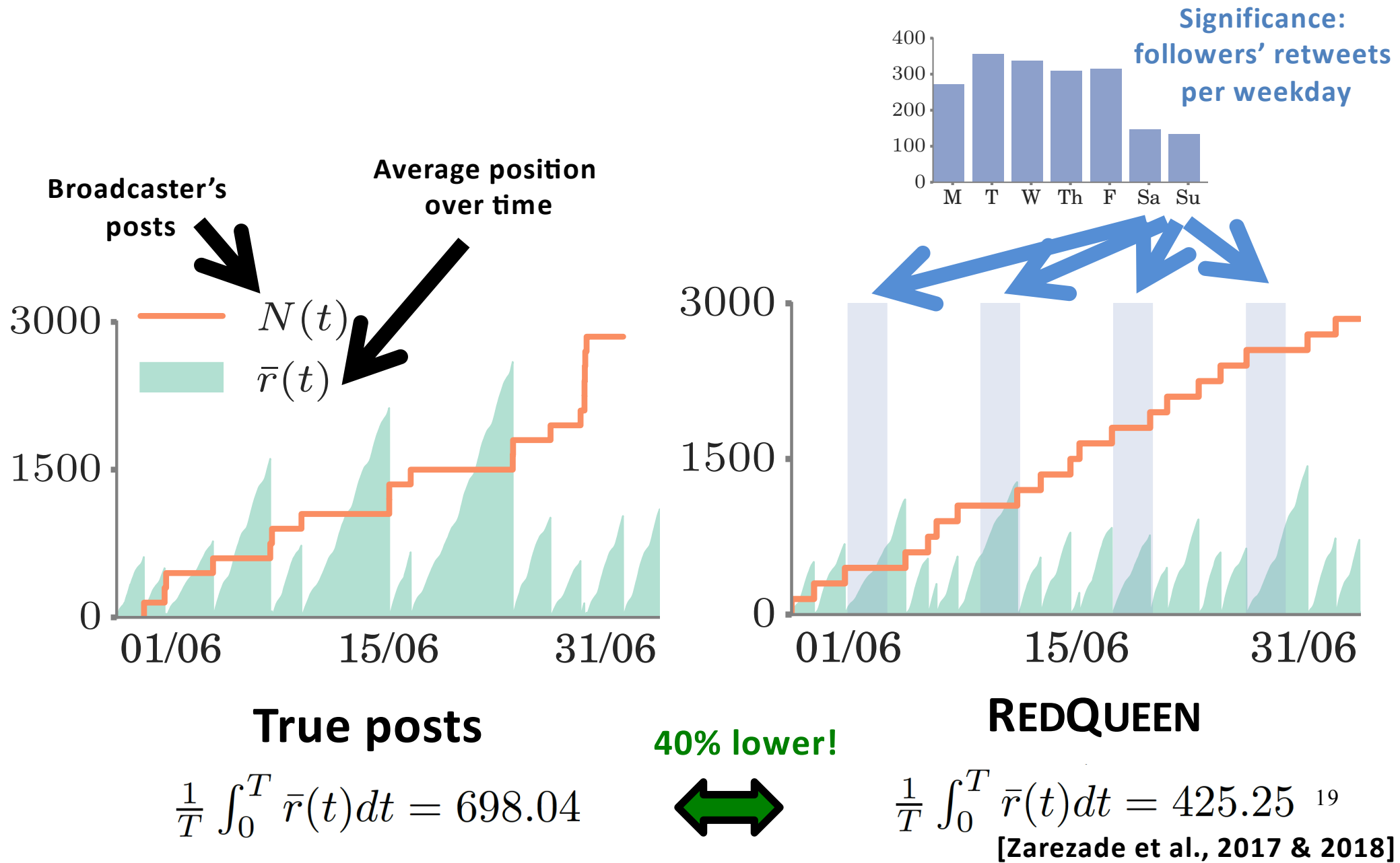
Consider $s(t) = s \Rightarrow u^*(t) = (s/q)^{1/2} r(t)$

How do we sample the next time?



It only requires sampling $M(t_f)$ times!

Example: a broadcaster in Twitter



RL & Control

1. Marked TPP: a new setting
2. Stochastic optimal control
- 3. Reinforcement learning**

Reinforcement learning of marked TPP

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



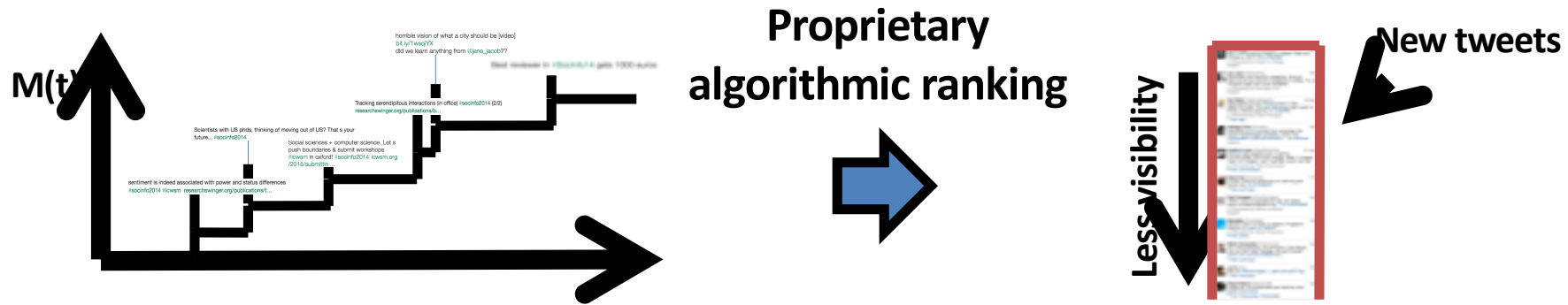
Next, details on one approach to the when/what to post problem with algorithmic ranking

oral

Similarly as with optimal control:

Policy is characterized by an intensity function!

Visibility dynamics are unknown



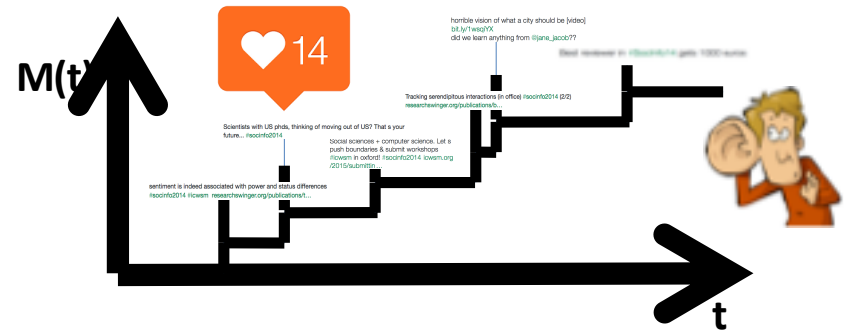
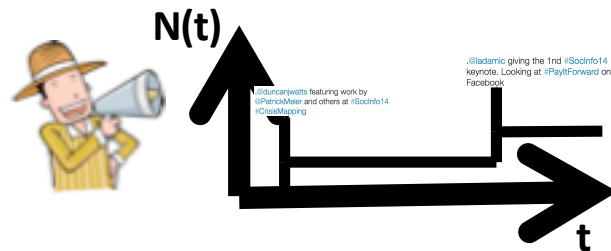
However, one may have access to quality metrics



Key idea:

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

Broadcasters and feedback



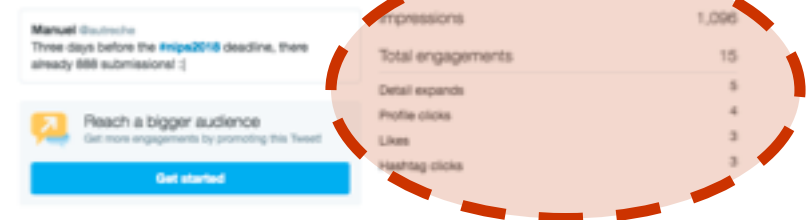
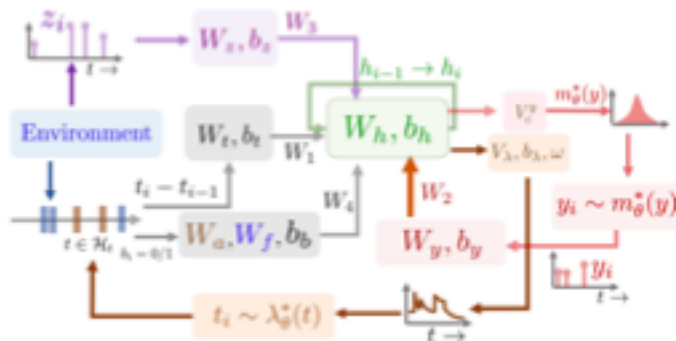
$$p_{\mathcal{A};\theta}^* = (\lambda_{\theta}^*, m_{\theta}^*)$$

Policy ↑ Intensity ↑ Mark distribution ↙

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

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

Parametrized using RNNs



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

Policy gradient

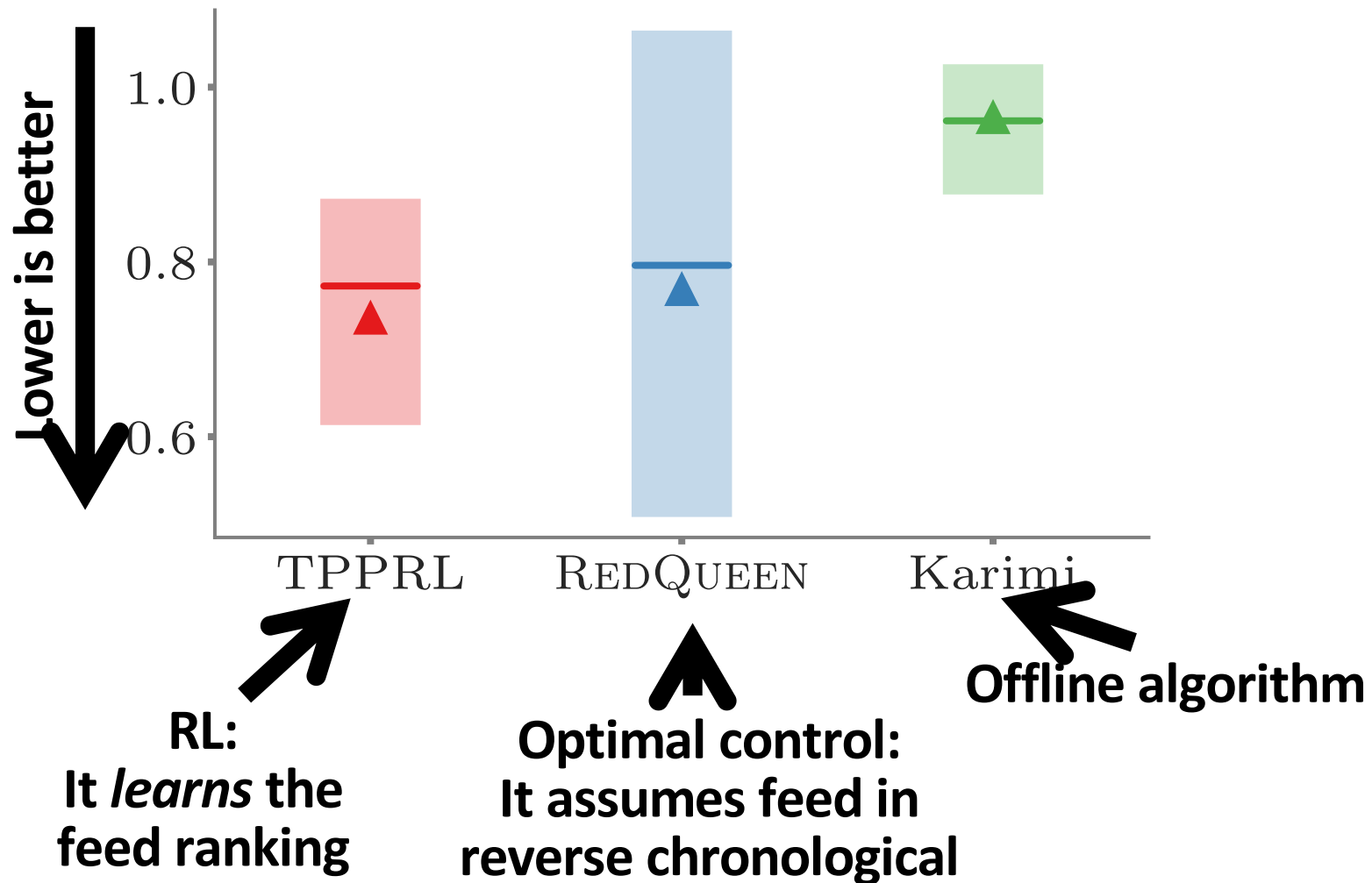
We aim to maximize the average reward in a time window $[0, T]$:

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

It can be shown that the reinforce trick is valid, i.e., we can compute the gradient and use SGD:

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

Example: 100 broadcasters in Twitter



Many thanks!

TEMPORAL POINT PROCESSES (TPPs): INTRO

1. Intensity function
2. Basic building blocks
3. Superposition
4. Marks and SDEs with jumps

MODELS & INFERENCE

1. Modeling event sequences
2. Clustering event sequences
3. Capturing complex dynamics
4. Causal reasoning on event sequences

RL & CONTROL

1. Marked TPPs: a new setting
2. Stochastic optimal control
3. Reinforcement learning