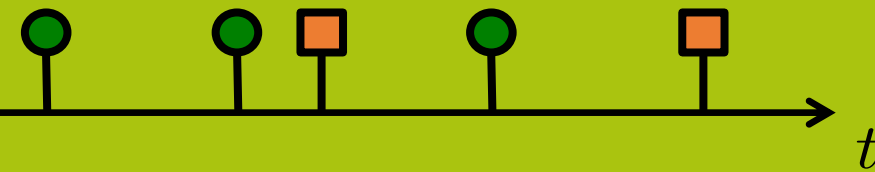# Learning with Temporal Point Processes: Models and Inference

**Alexey Zaytsev**
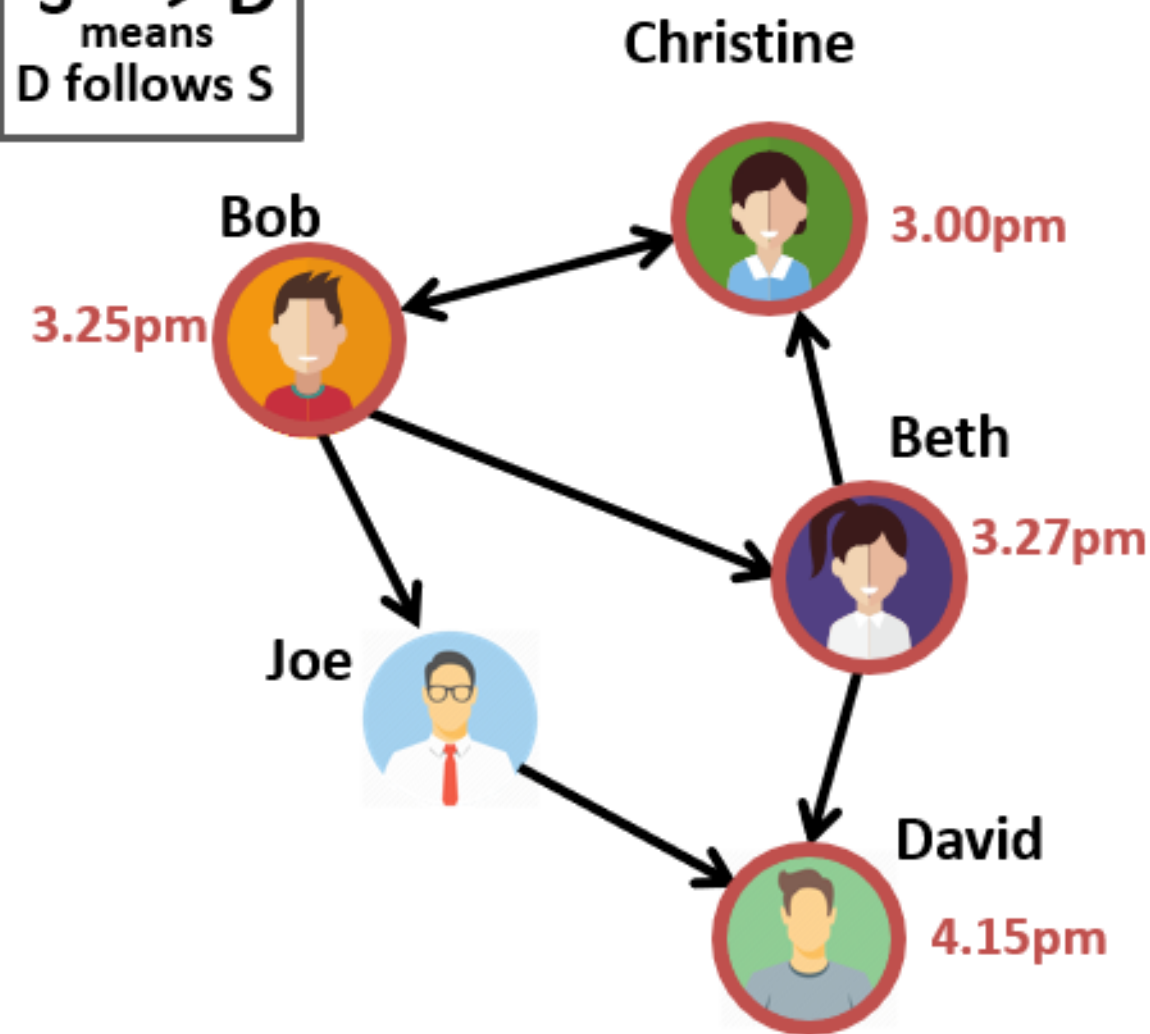
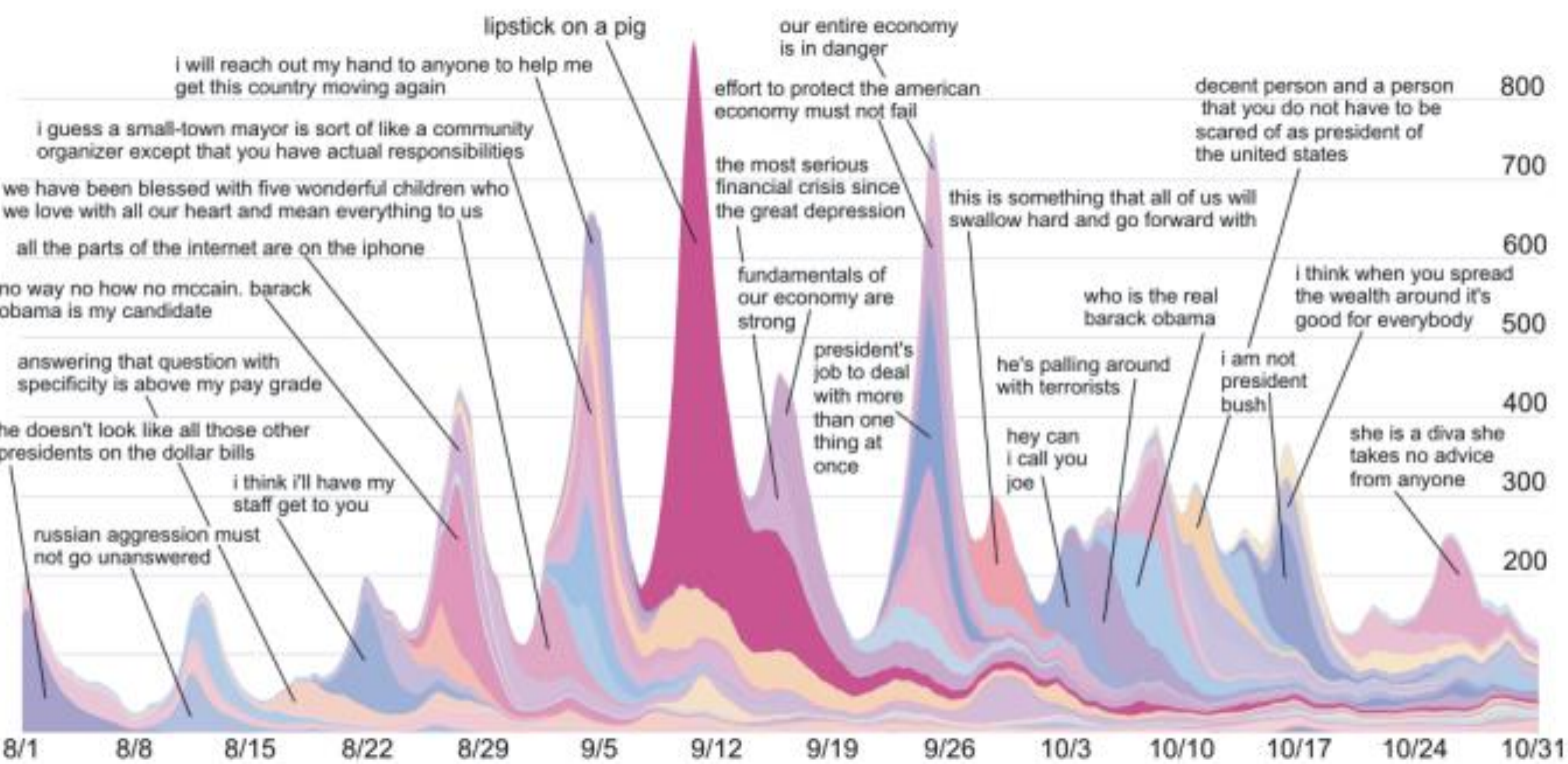Skoltech

# 1. Modeling event sequences

# Event sequences as cascades



S ⟶ D means D follows S
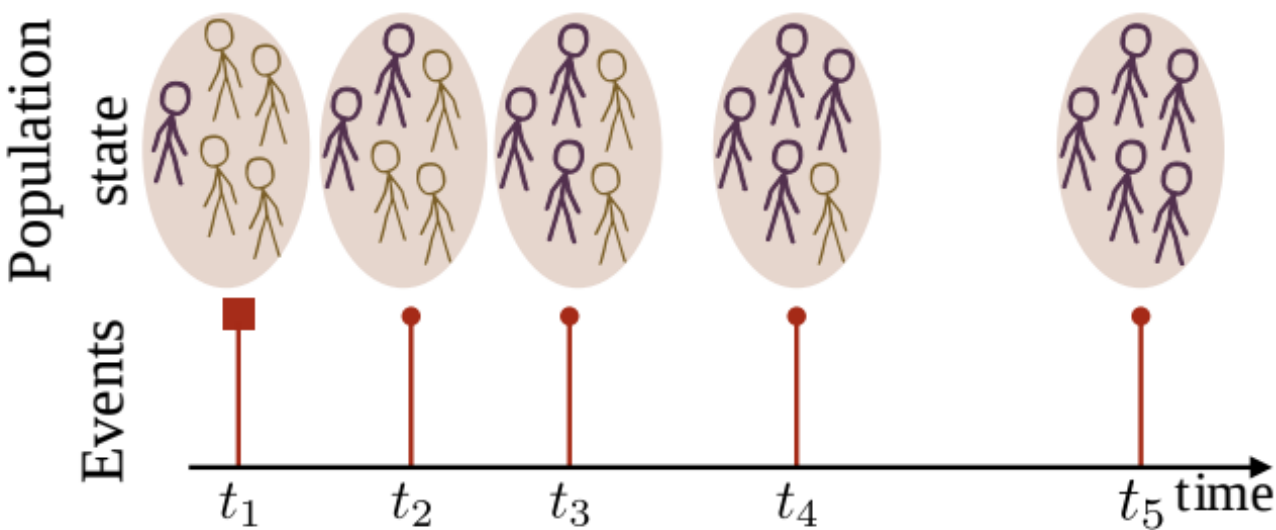
Christine — 3.00pm
Bob — 3.25pm
Beth — 3.27pm
Joe
David — 4.15pm

**Information Diffusion**



lipstick on a pig
our entire economy is in danger
i will reach out my hand to anyone to help me get this country moving again
effort to protect the american economy must not fail
decent person and a person that you do not have to be scared of as president of the united states
i guess a small-town mayor is sort of like a community organizer except that you have actual responsibilities
the most serious financial crisis since the great depression
this is something that all of us will swallow hard and go forward with
we have been blessed with five wonderful children who we love with all our heart and mean everything to us
all the parts of the internet are on the iphone
fundamentals of our economy are strong
i think when you spread the wealth around it's good for everybody
no way no how no mccain. barack obama is my candidate
president's job to deal with more than one thing at once
who is the real barack obama
answering that question with specificity is above my pay grade
he's palling around with terrorists
i am not president bush
he doesn't look like all those other presidents on the dollar bills
hey can i call you joe
she is a diva she takes no advice from anyone
i think i'll have my staff get to you
russian aggression must not go unanswered

8/1  8/8  8/15  8/22  8/29  9/5  9/12  9/19  9/26  10/3  10/10  10/17  10/24  10/31

800 700 600 500 400 300 200

[Leskovec et al., 2009]

**Disease Diffusion**



Population state

Events

$t_1$  $t_2$  $t_3$  $t_4$  $t_5$ time

[Rizoiu et al., 2018]



Event: $(t_i, u_i)$

Time    User

# An example: idea adoption



S →(means) D
D follows S

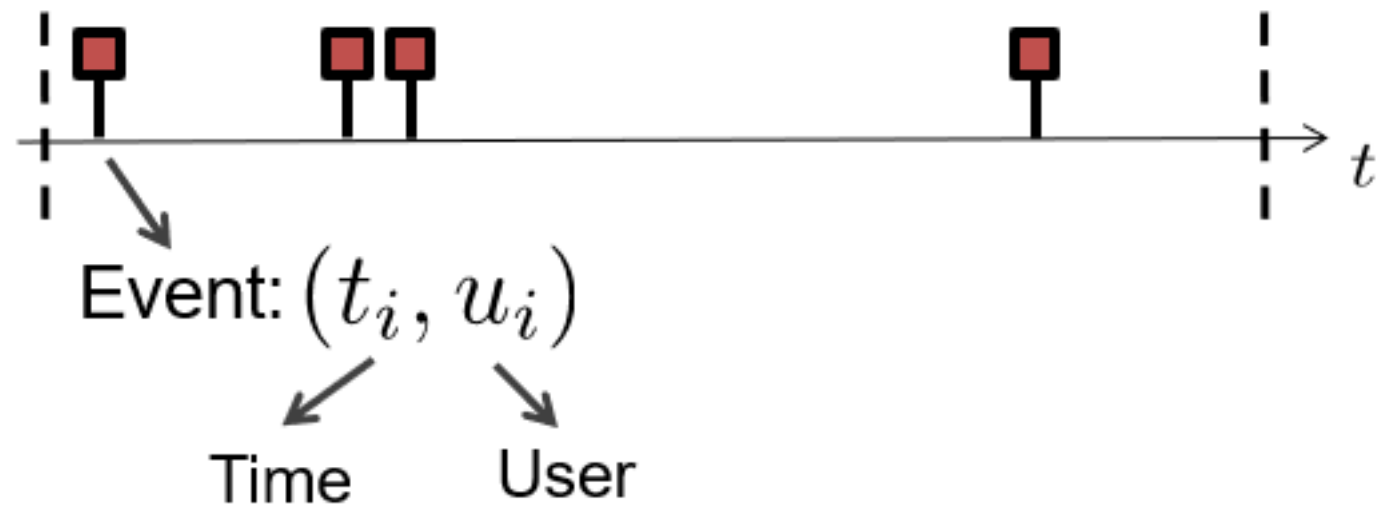Christine 3.00pm
Bob 3.25pm
Beth 3.27pm
Joe
David 4.15pm

Event: $(t_i, u_i)$
Time   User

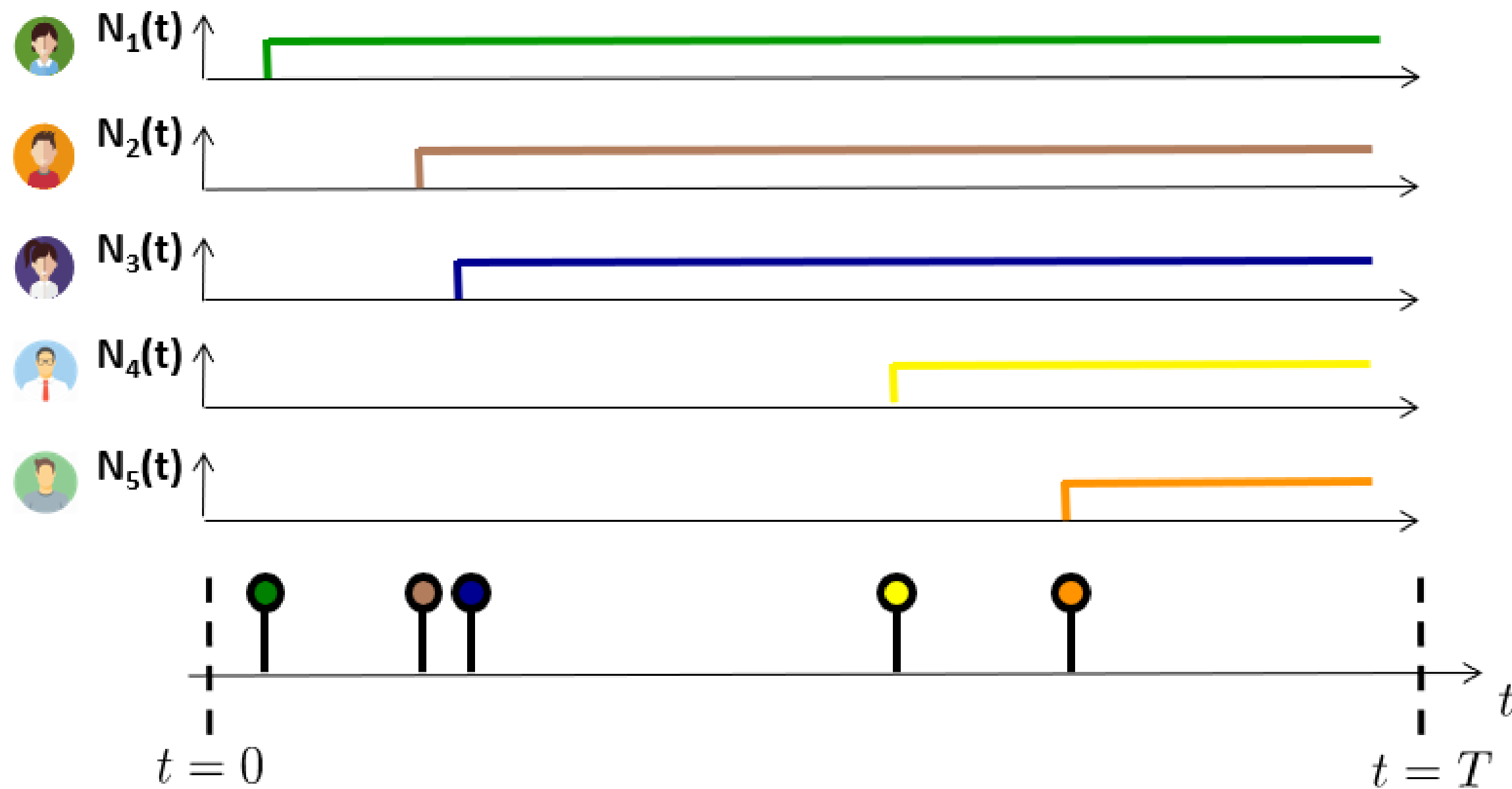They can have an impact
in the off-line world
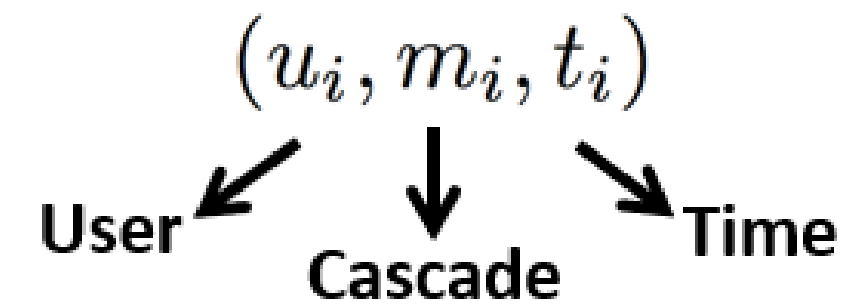


**theguardian**
Click and elect: how fake news helped
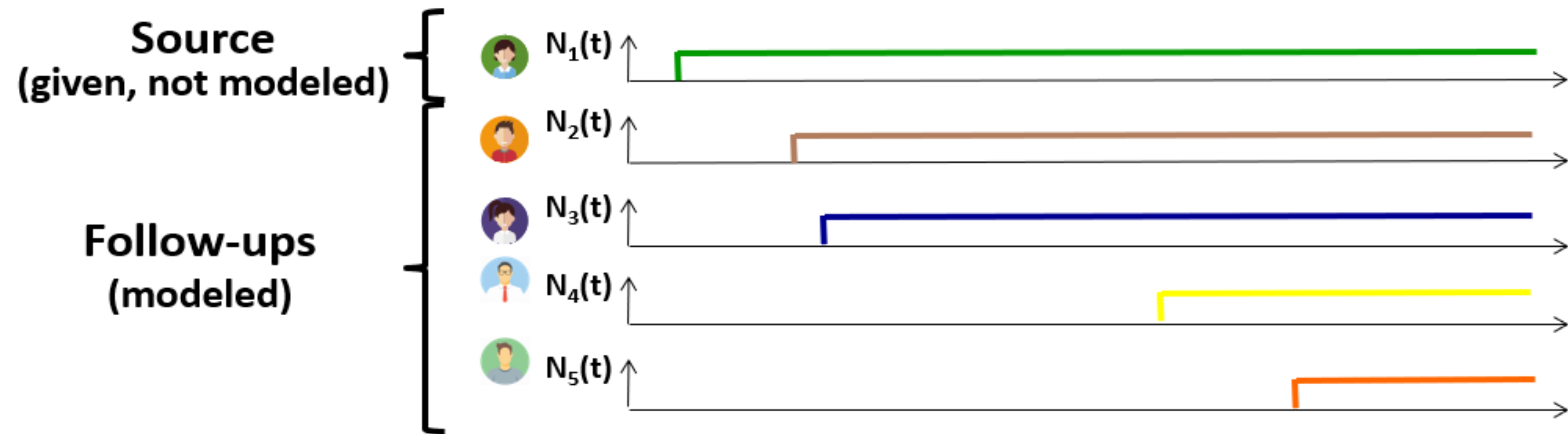Donald Trump win a real election

# Infection cascade representation

We represent an infection cascade using
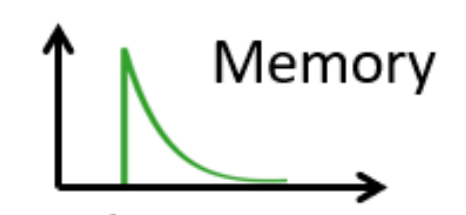**terminating temporal point processes**:



Infection event:

$$(u_i, m_i, t_i)$$

User    Cascade    Time

# Infection intensity

**Source**
(given, not modeled)

$N_1(t)$

$N_2(t)$

**Follow-ups**
(modeled)

$N_3(t)$

$N_4(t)$

$N_5(t)$

Memory

$$\lambda_u^*(t) = \underbrace{(1 - N_u(t))}_{\substack{\text{Users get} \\ \text{infected only} \\ \text{once}}} \underbrace{\sum_{v \in [m]} b_{vu}}_{\substack{\text{Influence from} \\ \text{user v on user u}}} \underbrace{\sum_{e_i \in \mathcal{H}_v(t)} \kappa(t - t_i)}_{\substack{\text{Previous} \\ \text{infections of user v}}}$$

[Gomez-Rodriguez et al., ICML 2011]

# Model inference from multiple cascades

**Conditional intensities** ⟷ **Diffusion log-likelihood**

$$\lambda_u^*(t)$$

$$\mathcal{L} = \sum_{u=1}^{n} \log \lambda_u^*(t_u) - \int_0^T \lambda_u^*(\tau)\,d\tau$$

**Maximum likelihood approach to find model parameters!**

**Sum up log-likelihoods of multiple cascades!**

**Theorem.** For any choice of parametric memory, the **maximum likelihood** problem is **convex in B**.

# Dynamic influence

In some cases, influence change over time:



**#greece retweets**

Propagation over networks with variable influence

Properties are similar to static influence



**0**

**T**

[Gomez-Rodriguez et al., WSDM 2013]

Skoltech

# Recurrent events: beyond cascades

**Up to this point,** each user is only infected once, and event sequences can be seen as cascades.

**In general, users perform recurrent events over time.** E.g., people repeatedly express their opinion online:

CNN
How social media is revolutionizing debates

The New York Times
*Social Media Are Giving a Voice to Taste Buds*

TC
**Twitter Unveils A New Set Of Brand-Centric Analytics**

The New York Times
*Campaigns Use Social Media to Lure Younger Voters*

Skoltech

# Recurrent events representation

We represent messages using **nonterminating temporal point processes**:



Recurrent event:

$$(u_i, t_i)$$

User ↙    ↘ Time

[Farajtabar et al., NIPS 2014]

# Recurrent events intensity

$N_1(t)$

$N_2(t)$

$N_3(t)$

$N_4(t)$

$N_5(t)$

**Cascade sources!**

$$\lambda_u^*(t) \;=\; \mu_u \;+\; \sum_{v \in [m]} b_{vu} \sum_{e_i \in \mathcal{H}_v(t)} \kappa(t - t_i)$$

Memory

Hawkes process

User's intensity

Events on her own initiative

Influence from user v on user u

Previous messages by user v

[De et al., NIPS 2016]

# 2. Clustering event sequences

Skoltech

# Event sequences

So far, we have assumed the cascade (topic, meme, etc.) that each event belongs to was known.

Often, the cluster (topic, meme, etc.) that each event in a sequence belongs to is not known:



**BBC News (World)** ✔ @BBCWorld · 4m
Turkey election: Erdogan win ushers in new presidential era

Politics

**BBC News (World)** ✔ @BBCWorld · 46m
Dublin church: Seven injured as car hits pedestrians

**BBC News (World)** ✔ @BBCWorld · 2h
Nigerian music star D'banj's son 'drowns at home'

Music

**BBC News (World)** ✔ @BBCWorld · 2h
Turkey election: Country's heart split over Erdogan victory

Politics

# Clustering event sequences

Assume the event <u>cluster to be hidden</u> and aim to automatically <u>learn the cluster assignments</u> from the data:



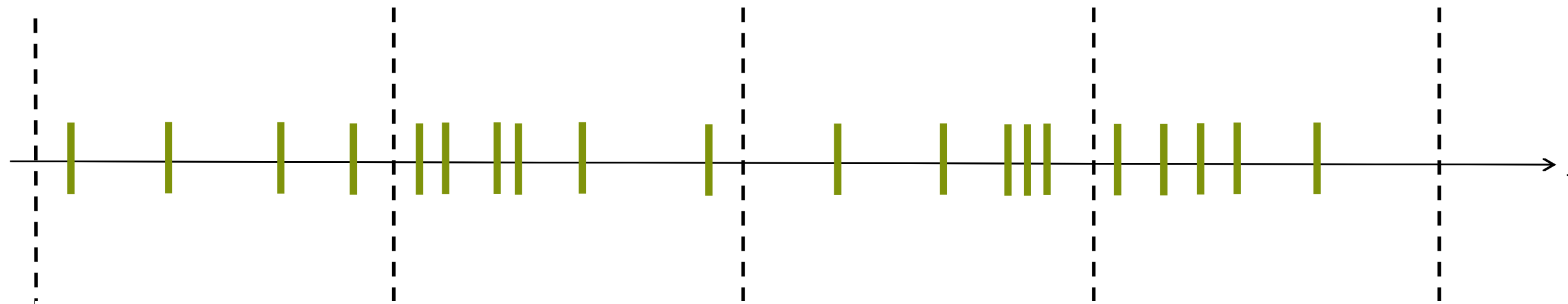<u>Bayesian methods</u> to cluster event sequences in the context of:



**Learning**



**Online News**

| Method | DMHP |
|---|---|
| ICU Patient | 0.3778 |
| IPTV User | 0.2004 |

**Health care**

[Du et al., 2015; Mavroforakis et al., 2017; Xu & Zha, 2017]
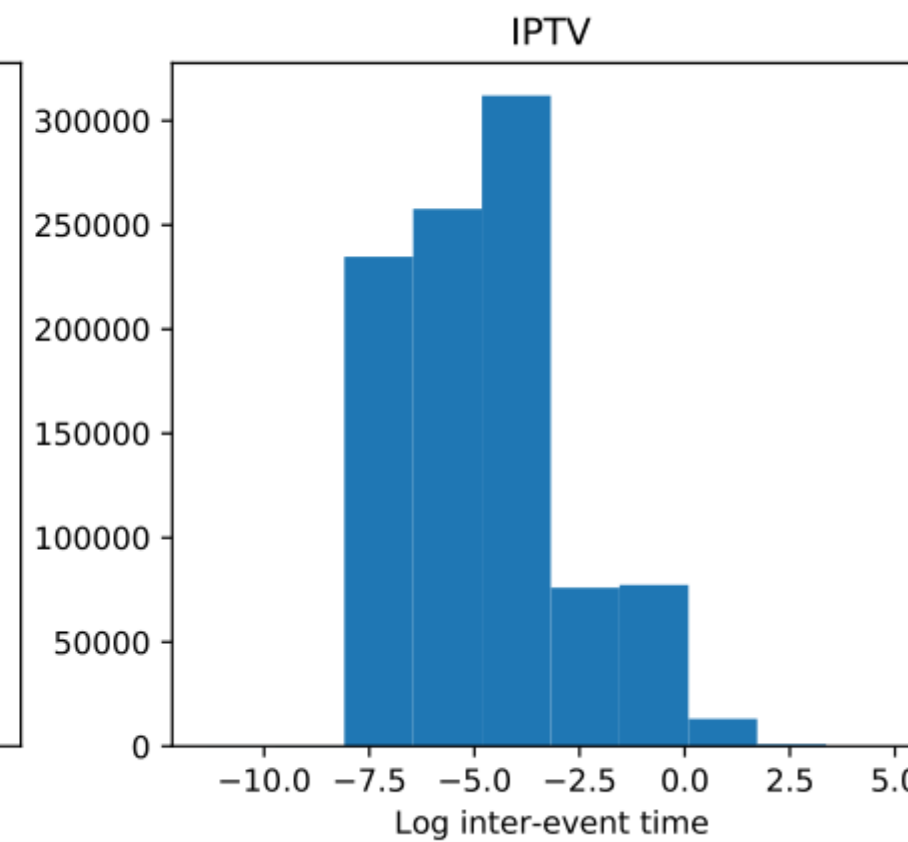
# Clustering event sequences

Assume the event <u>cluster to be hidden</u> and aim to automatically <u>learn the cluster assignments</u> from the data:



**Financial transactions**          **Job changes from Linkedin**          **TV programs, IPTV**

[Zhuzhel et al., 2021]

# Hierarchical Dirichlet Hawkes process

**1st year computer science student**

- Introduction to programming
- Discrete math
- Project presentation



[Mavroforakis et al., WWW 2017]

# Events representation

We represent the events using <u>marked temporal point processes</u>:



$N_{u,\ell}(t)$

Task

Task

$t = 0$

$t = T$
$N_{u,\ell}(T) = 9$

$t$

Event: $(t_n, p_n, q_n)$

Time

Cluster
(hidden)

Content

[Mavroforakis et al., WWW 2017]

Skoltech

# Cluster intensity



New cascade rate    Cluster popularity    Memory

$$\lambda^*_{u,\ell}(t) = \mu_u \pi_\ell + \sum_{j:t_j \in \mathcal{H}_{u,\ell}(t)} k_{\theta_\ell}(t - t_j)$$

Intensity or rate (events / hour)    Own initiative    Follow-up

Hawkes process

$N_{u,\ell}(t)$

Task    Task

Event: $(t_n, p_n, q_n)$

$t = 0$    $t = T$
$N_{u,\ell}(T) = 9$

Time    Content

Cluster (hidden)

Skoltech

[Mavroforakis et al., WWW 2017]

# User events intensity

Users adopt more than one cluster:



A user's learning events as a multidimensional Hawkes**:**

$$\underset{\substack{\nearrow \\ Time}}{(t_n,} \underset{\substack{\nwarrow \\ cluster}}{p_n)} \sim Hawkes \begin{pmatrix} \lambda^*_{u,1}(t) \\ \vdots \\ \lambda^*_{u,\infty}(t) \end{pmatrix}$$

$$Content \rightarrow q_n = \boldsymbol{\omega} \quad \omega_j \sim Multinomial(\boldsymbol{\theta}_p)$$

Skoltech

[Mavroforakis et al., WWW 2017]

# People share same clusters

Different users adopt same clusters



Cluster distribution from a <u>Dirichlet process</u>:

- Infinite # of clusters.
- Shared parameters across users.

Efficient model inference using Sequential Monte-Carlo!

Skoltech

[Mavroforakis et al., WWW 2017]
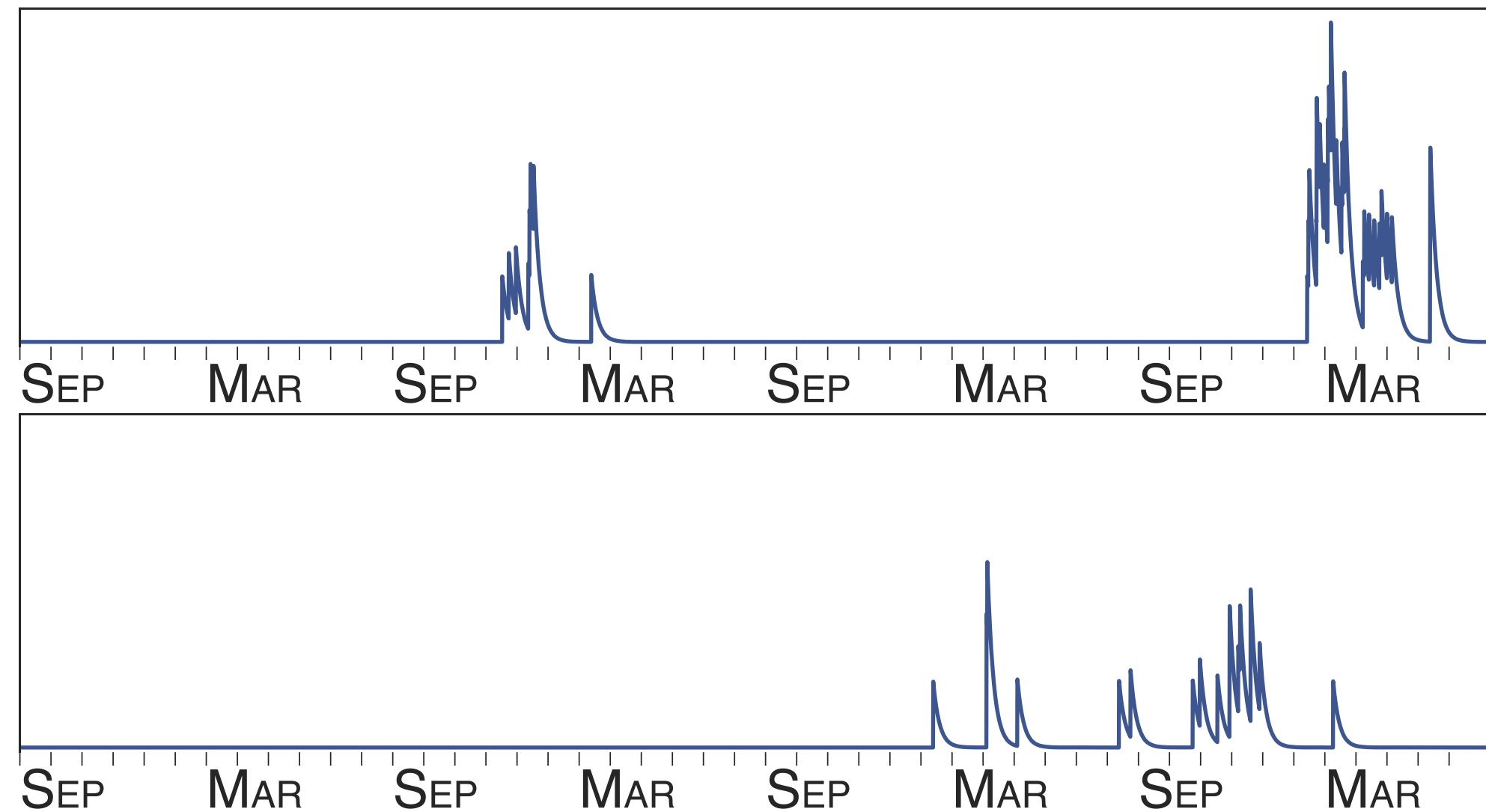
# Learning cluster (I): Version Control

Content

Intensities



Version control tasks tend to be specific,
quickly solved after performing few questions
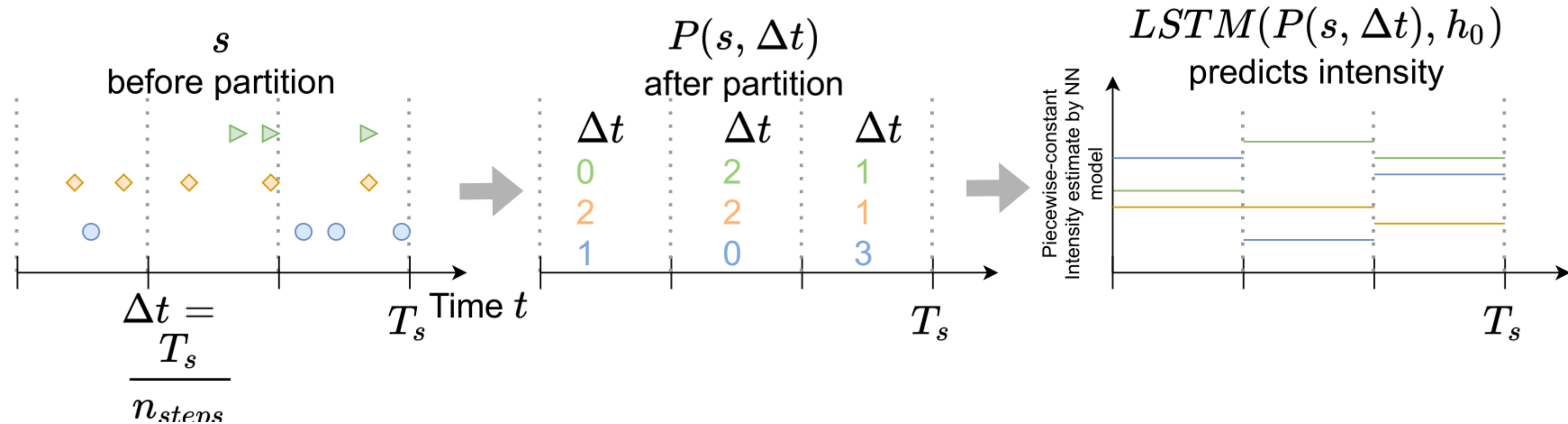
# Learning cluster (II): Machine learning

Content

Intensities



Machine learning tasks tend to be more
complex and require asking more questions

[Mavroforakis et al., WWW 2017]

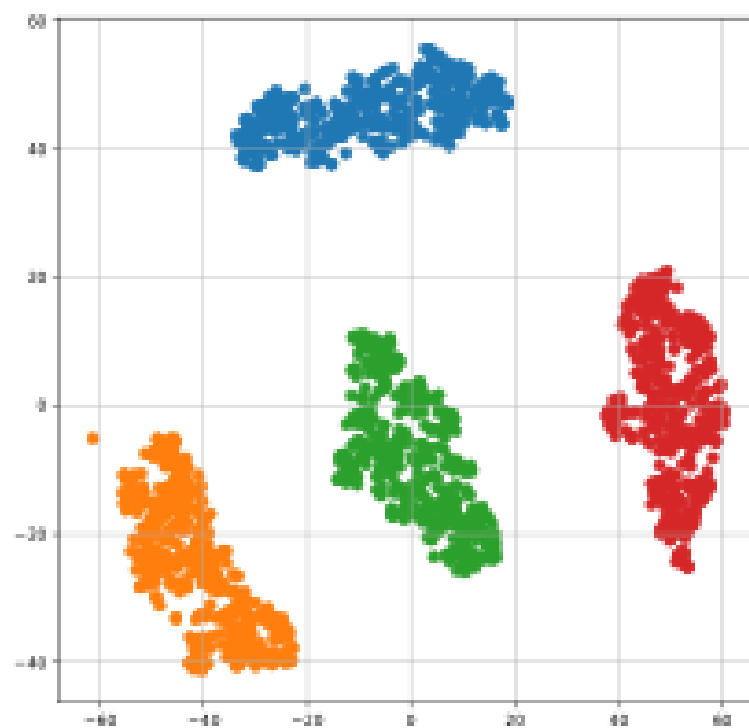# COHORTNEY for events clustering



Analytical likelihood

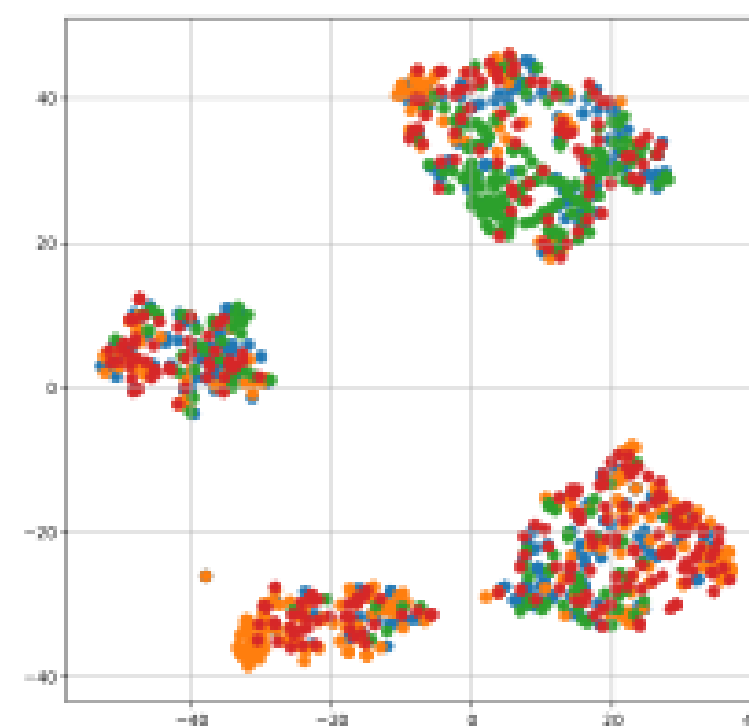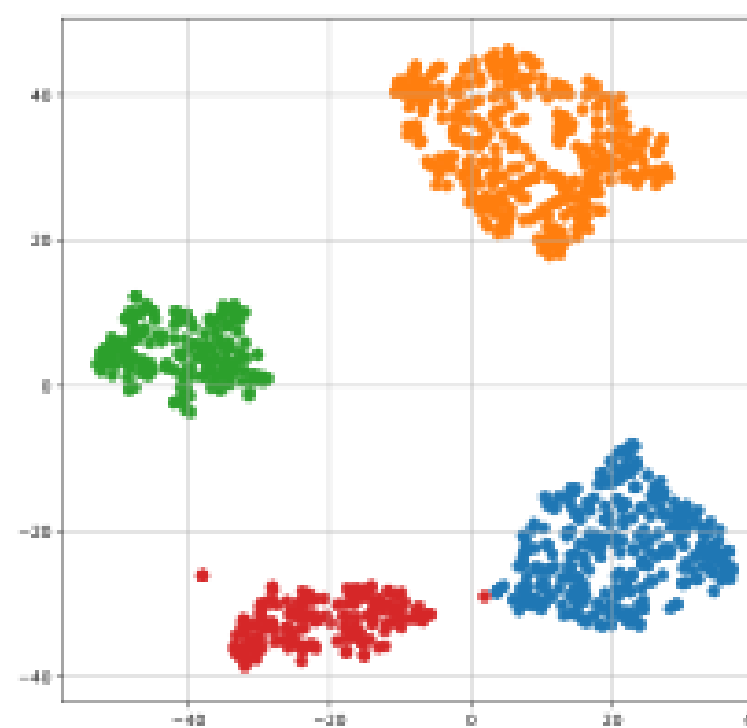EM algorithm for selection of parameters and labeling

[Zhuzhel et al., 2021]

# COHORTNEY for events clustering



(a)

(b)

(c)

(d)

True (a) and learned (b) clusters
for synthetic data

True (c) and learned (d) clusters
for real AGE data

Skoltech

[Zhuzhel et al., 2021]

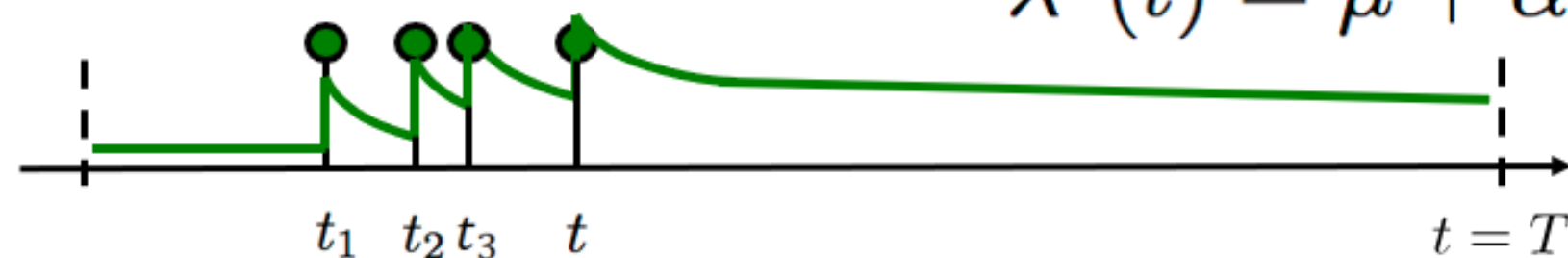# 3. Capturing complex dynamics

Skoltech

# Towards real-world temporal dynamics

Up to now, we have focused on simple temporal
dynamics (and intensity functions):

$$\lambda^*(t) = \mu$$

$$\lambda^*(t) = \sum_j \alpha_j\, k(t - t_j)$$

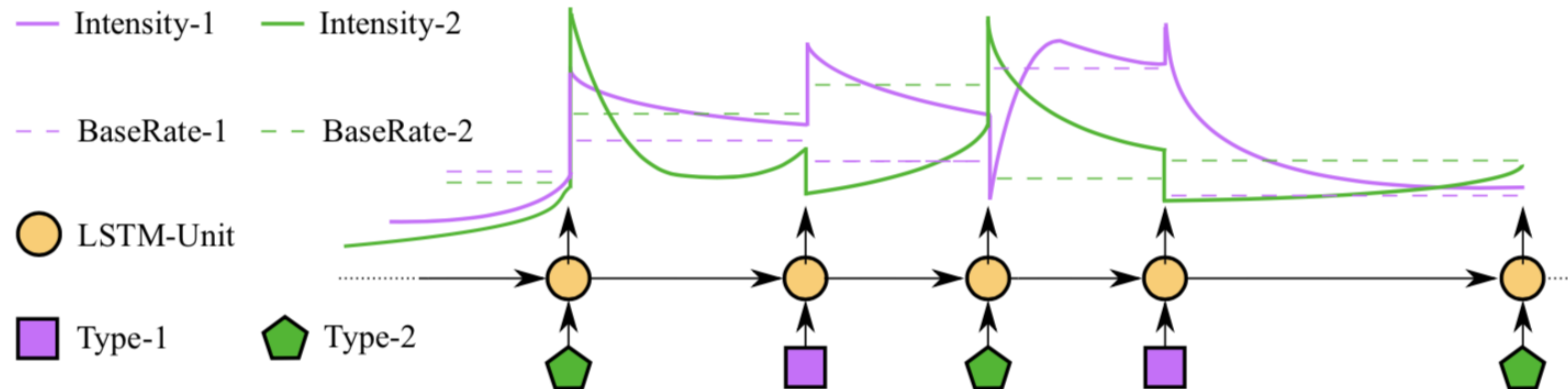$$\lambda^*(t) = \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i)$$

Recent works make use of RNNs to capture
more complex dynamics

[Du et al., 2016; Dai et al., 2016; Mei & Eisner, 2017; Jing & Smola,
2017; Trivedi et al., 2017; Xiao et al., 2017a; 2018]

# Neural Hawkes process

1) History effect does not need to be additive

2) Allows for complex memory effects

   such as delays
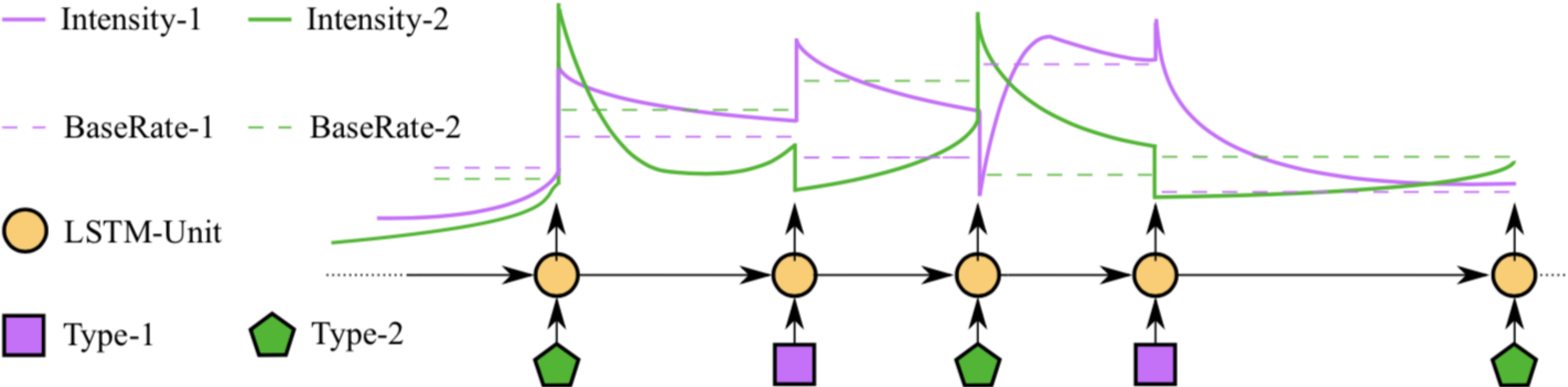


[Mei & Eisner, NIPS 2017]

# Neural Hawkes process

$$\lambda_u(t) = f_u(\mathbf{w}_u^\top \mathbf{h}(t)) \qquad \mathbf{h}(t) = \overbrace{\mathrm{RNN}(\mathcal{H}(t))}^{\text{Memory}}$$

Excitation & inhibition

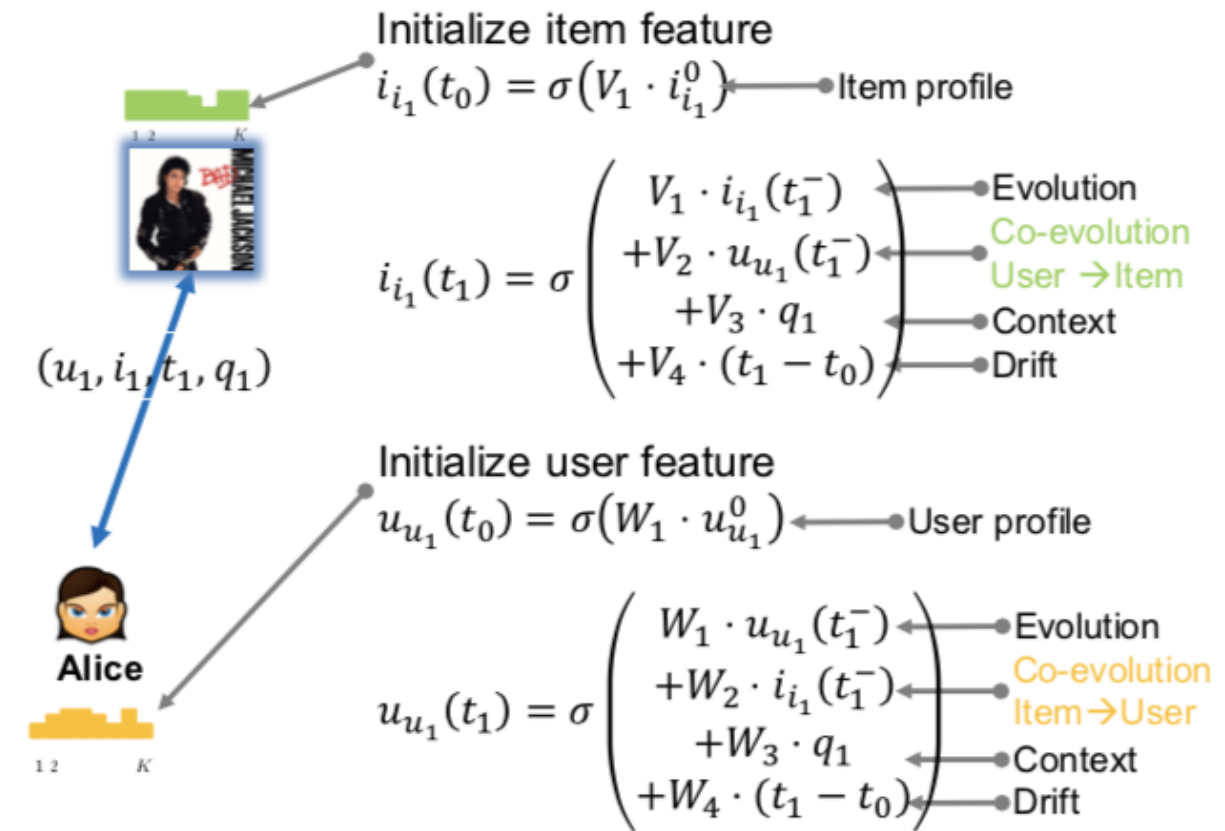Parametric learning using stochastic gradient descent



- Intensity-1
- Intensity-2
- BaseRate-1
- BaseRate-2
- LSTM-Unit
- Type-1
- Type-2

[Mei & Eisner, NIPS 2017]
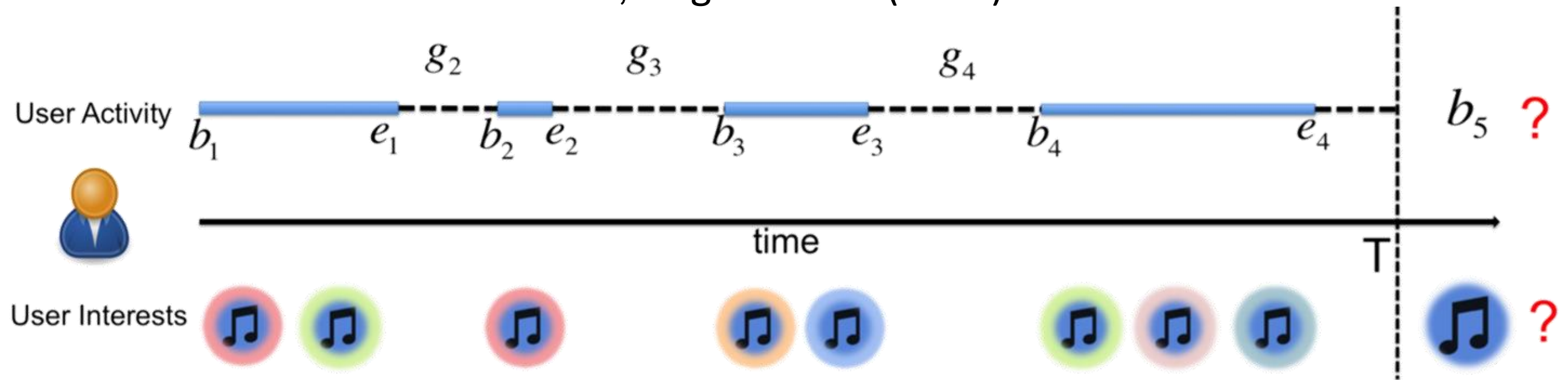
# Applications (I): Predictive Models

*Know-Evolve*, Trivedi et al. (2017)



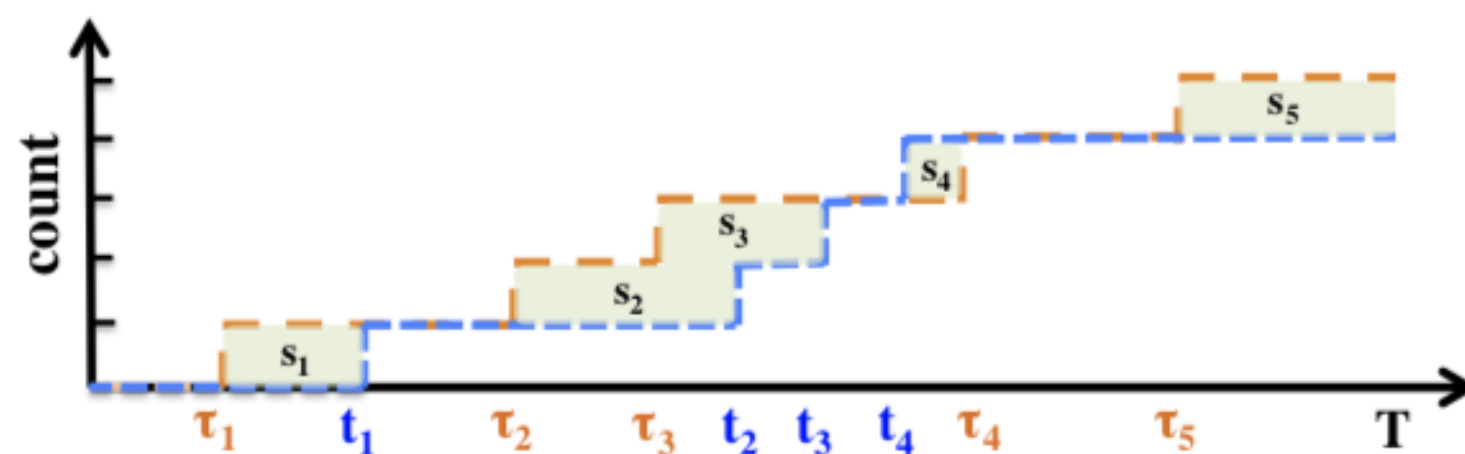*Coevolutionary Embedding*, Dai et al. (2017)



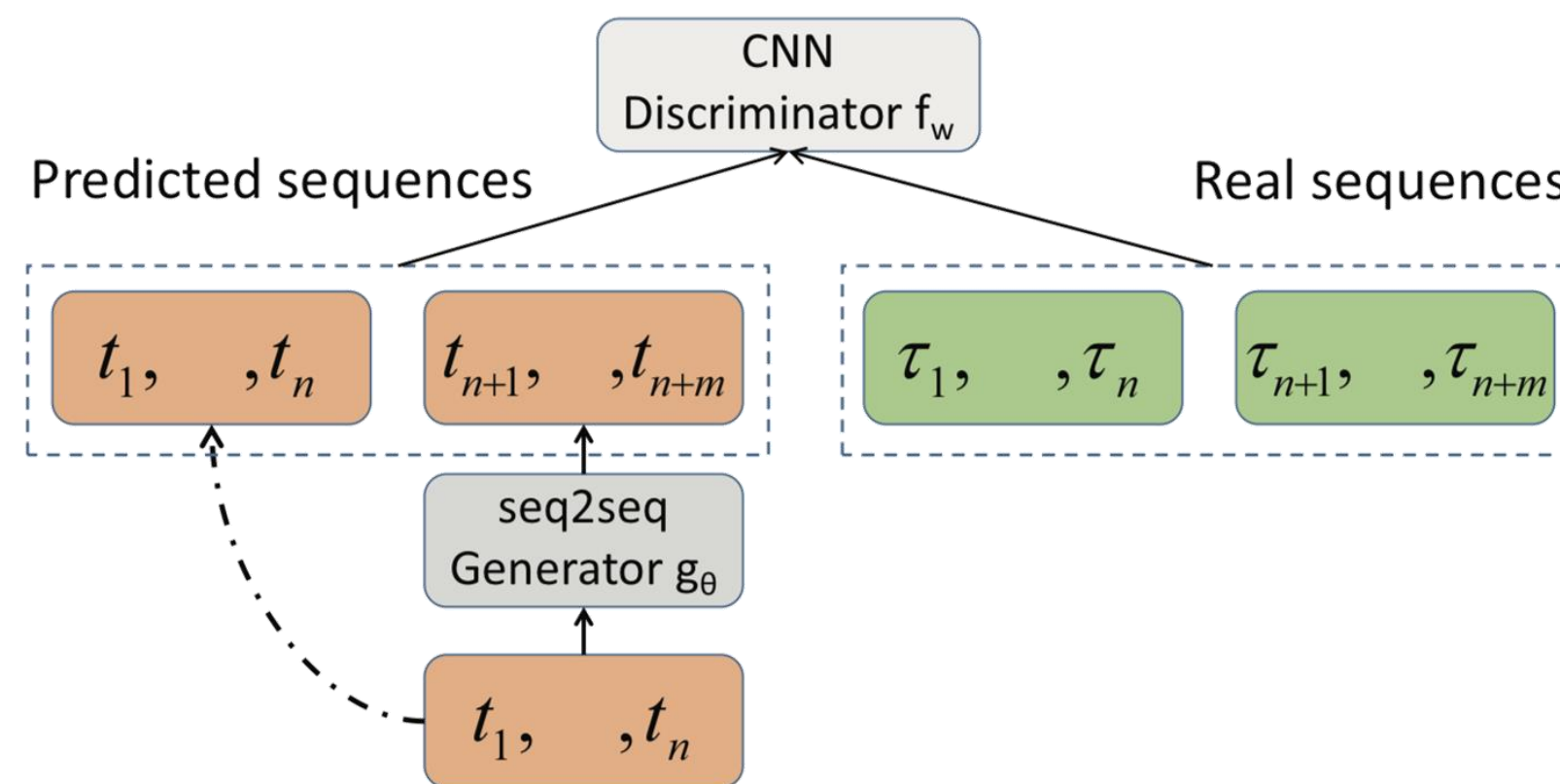*Neural Survival Recommender*, Jing & Smola (2017)

# Applications (II): Generative Models

Key idea: Intensity- and likelihood-free models



Wasserstein-Distance for
Temporal Point Processes

GAN architecture

[Xiao et al., 2017 & 2018]

# 4. Causal reasoning on event sequences

Skoltech

# Temporal point processes beyond prediction

So far, we have focused on models that improve predictions:
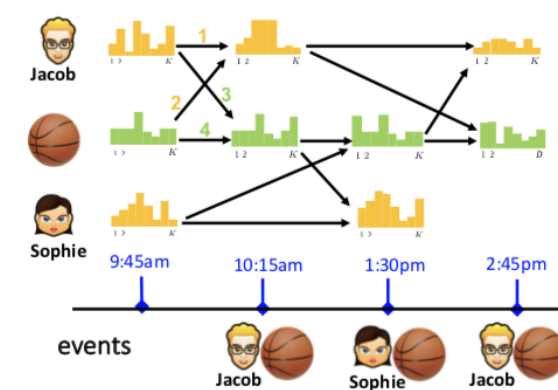
Link prediction

Community detection

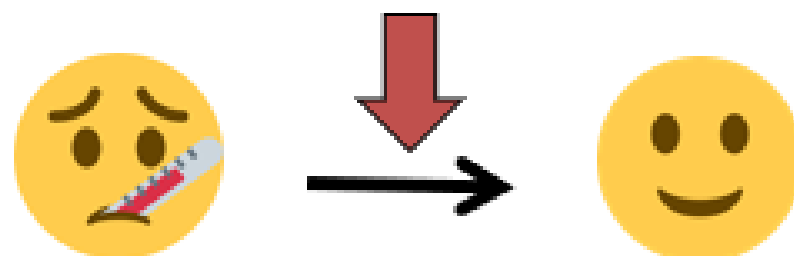Recommendations

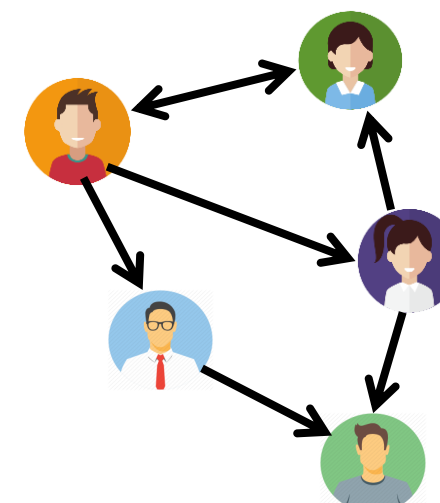[Trivedi et al., 2017]

[Xiao et al., 2017]

[Dai et al., 2017]

Recent works have focused on performing <u>causal inference using event sequences</u>:

**Treatment effect**

Granger causality graph

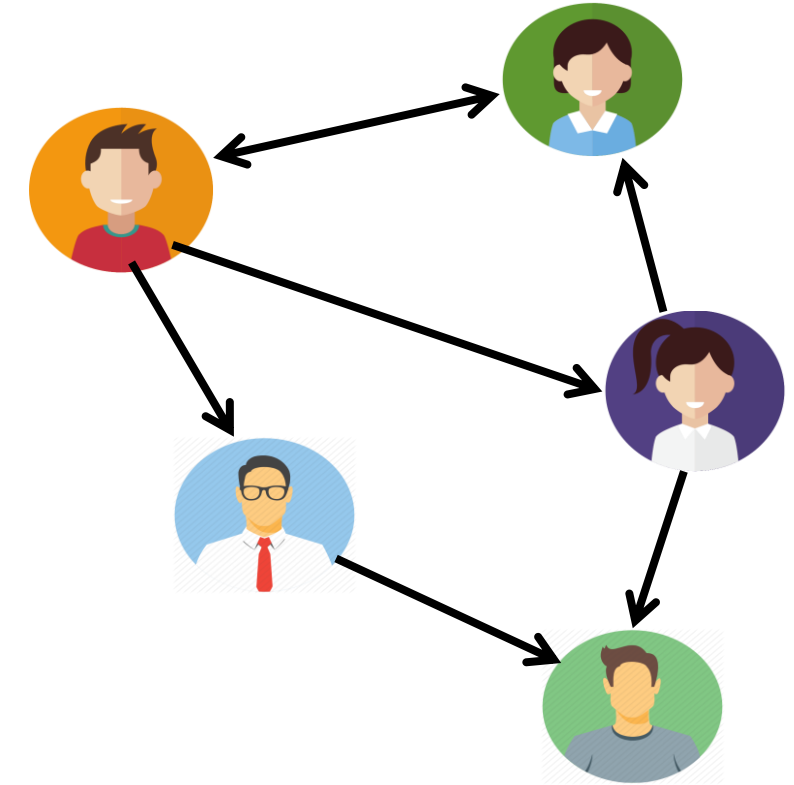[Xu et al., 2016; Achab et al., 2017; Kuśmierczyk & Gomez-Rodriguez, 2018]

# Uncovering Causality from Hawkes Processes

Multivariate Hawkes process:



$$N(t) = \sum_{u \in \mathcal{U}} N_u(t)$$

$$\lambda_u(t) = \mu_u + \sum_{v \in \mathcal{U}} \int_0^t \underbrace{k_{u,v}(t - t')dN_v(t')}$$

Effect of v's past events on u

**Granger causality:**

"X causes Y in the sense of Granger causality if forecasting future values of Y is more successful while taking X past values into account"
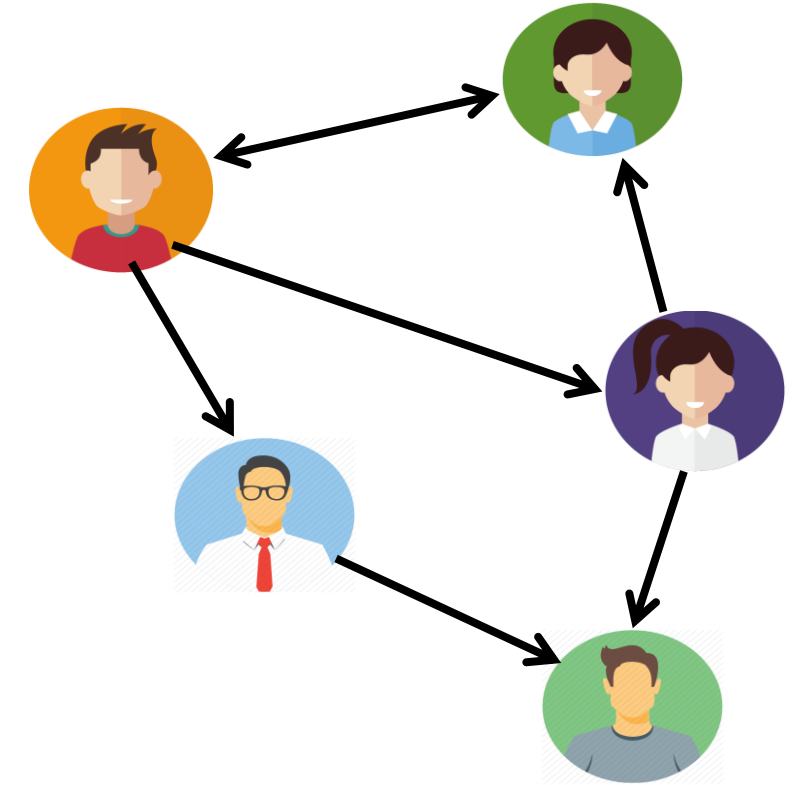
[Granger, 1969]

[Achab et al., ICML 2017]

# Uncovering Causality from Hawkes Processes

Multivariate Hawkes process:

$$N(t) = \sum_{u \in \mathcal{U}} N_u(t)$$

$$\lambda_u(t) = \mu_u + \sum_{v \in \mathcal{U}} \int_0^t \underbrace{k_{u,v}(t - t')dN_v(t')}$$

Effect of v's past events on u



**Granger causality on multivariate Hawkes processes:**

" $N_v(t)$ does not Ganger-cause $N_u(t)$ w.r.t. $N(t)$ if and only if

$$k_{u,v}(\tau) = 0 \text{ for } \tau \in \Re^+ \text{"}$$

[Eichler et al., 2016]
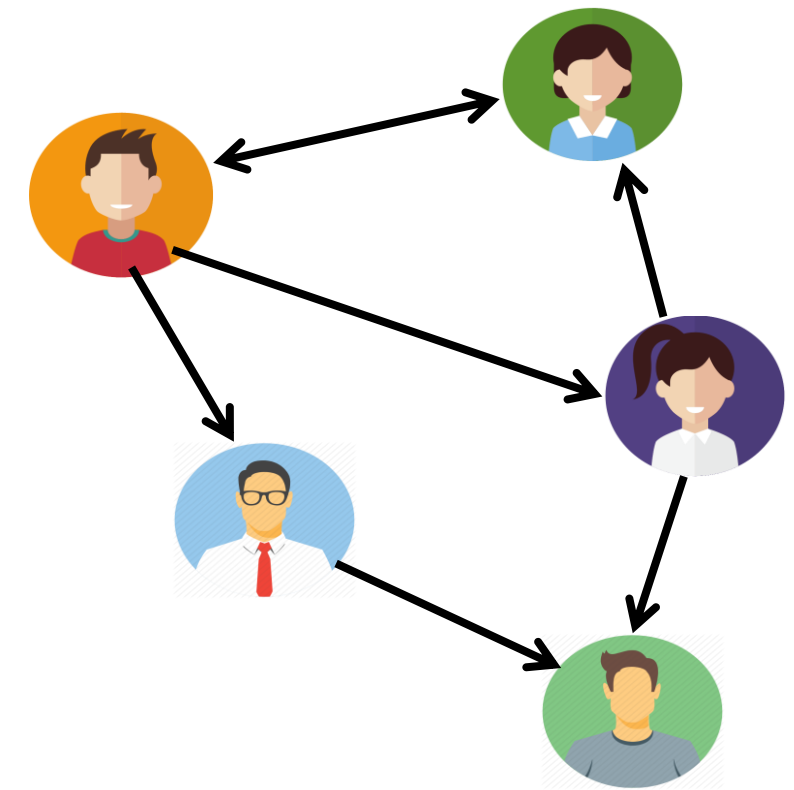
[Achab et al., ICML 2017]

# Uncovering Causality from Hawkes Processes

Goal is to estimate $G = [g_{uv}]$, where:

$$g_{uv} = \int_0^{+\infty} k_{u,v}(\tau) d\tau \geq 0 \text{ for all } u, v \in \mathcal{U}$$

Average total # of events of node *u* whose *direct* ancestor is an event by node *v*

Then, $G = [g_{uv}]$ quantifies the *direct causal relationship* between nodes.
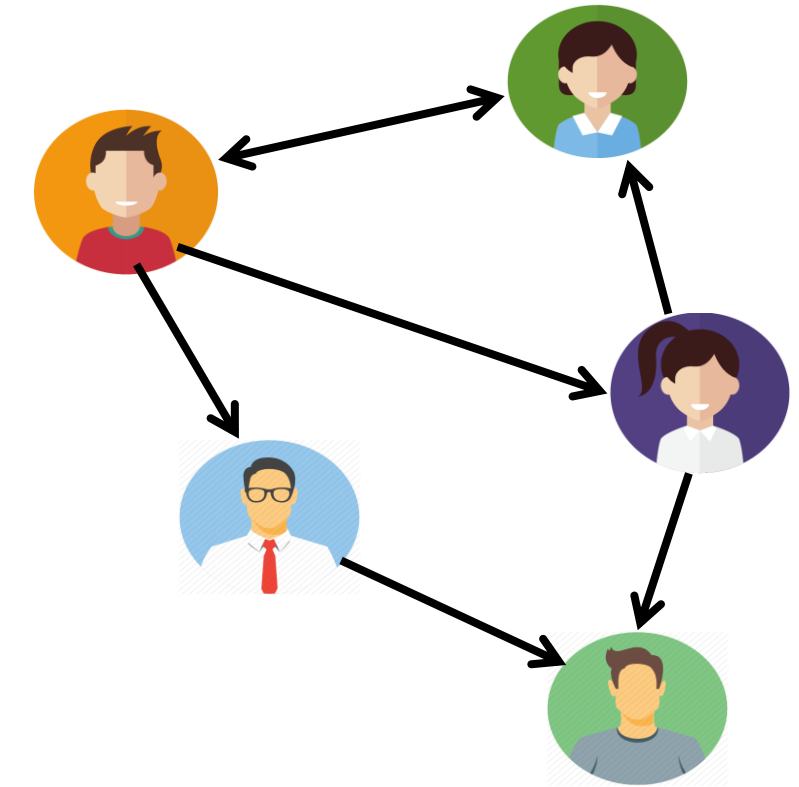


[Achab et al., ICML 2017]

# Uncovering Causality from Hawkes Processes

Goal is to estimate $G = [g_{uv}]$, where:

$$g_{uv} = \int_0^{+\infty} k_{u,v}(\tau)d\tau \geq 0 \text{ for all } u, v \in \mathcal{U}$$

Average total # of events of node *u* whose *direct* ancestor is an event by node *v*

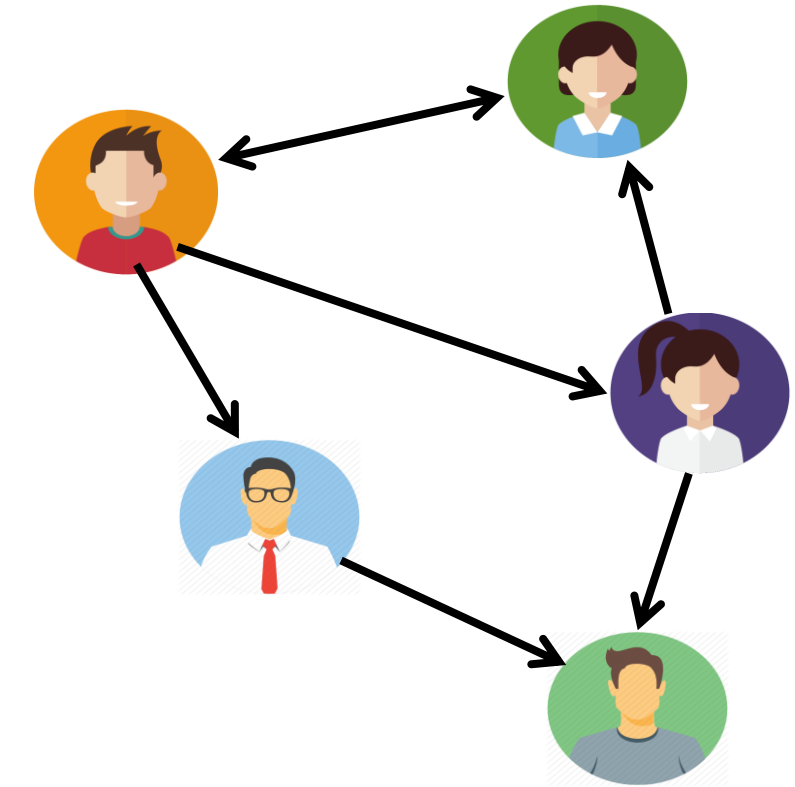Then, $G = [g_{uv}]$ quantifies the *direct causal relationship* between nodes.

Key idea: Estimate G using the cumulants dN(t) of the Hawkes process.

[Achab et al., ICML 2017]

# Uncovering Causality from Hawkes Processes

Goal is to estimate $G = [g_{uv}]$, where:

$$g_{uv} = \int_0^{+\infty} k_{u,v}(\tau)d\tau \geq 0 \text{ for all } u, v \in \mathcal{U}$$

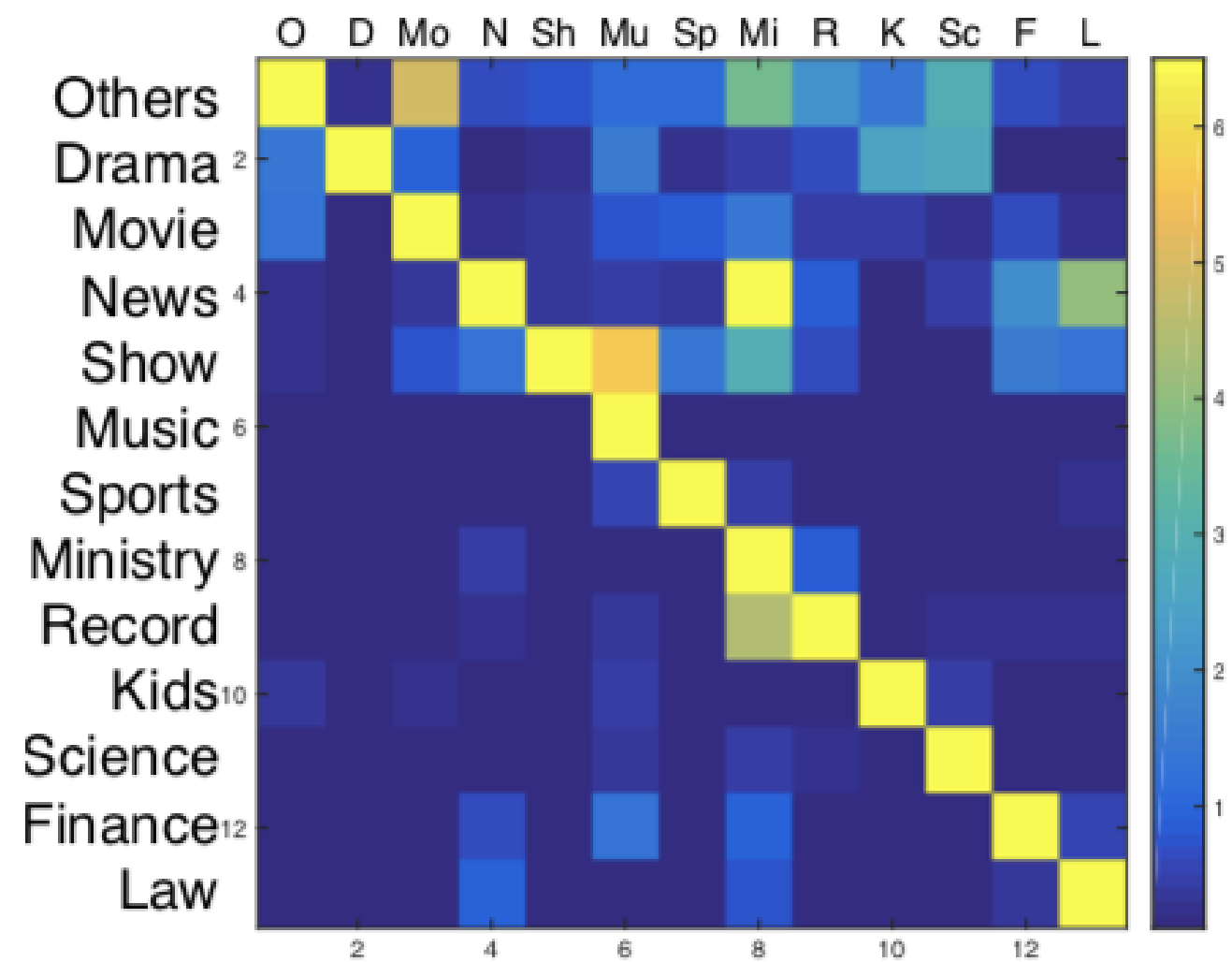Average total # of events of node *u* whose *direct* ancestor is an event by node *v*

Then, $G = [g_{uv}]$ quantifies the *direct causal relationship* between nodes.

Key idea: Estimate G using the cumulants dN(t) of the Hawkes process.

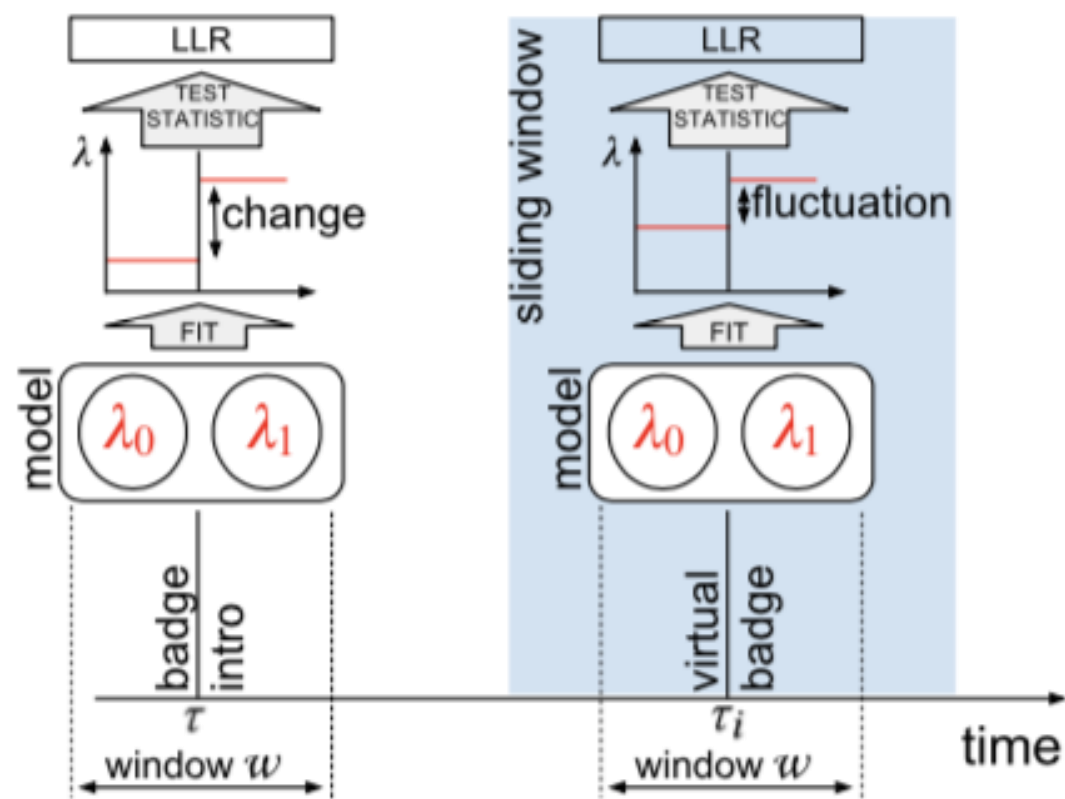Non parametric Hawkes cumulant estimation method with TensorFlow implementation

[Achab et al., ICML 2017]

Skoltech

# Causal reasoning: Applications

Infectivity matrix estimation

Effect of Badges



[Xu et al., 2016, ICML]

[Kuśmierczyk & Gomez-Rodriguez, 2018]

Tag wiki rank over time

Skoltech

# Models and inference

1. Modeling event sequences

2. Clustering event sequences

3. Capturing complex dynamics

4. Causal reasoning on event sequences

Skoltech