

# Neural Networks for point processes



**Alexey Zaytsev**

Assistant professor, Skoltech

22th of December

Based on ICML Tutorial, July 2018 by I. Valera & M.G. Rodriguez  
Slides/references: <http://learning.mpi-sws.org/tpp-icml18>



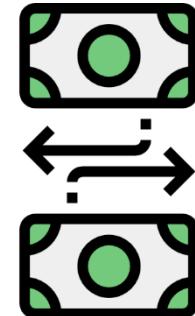
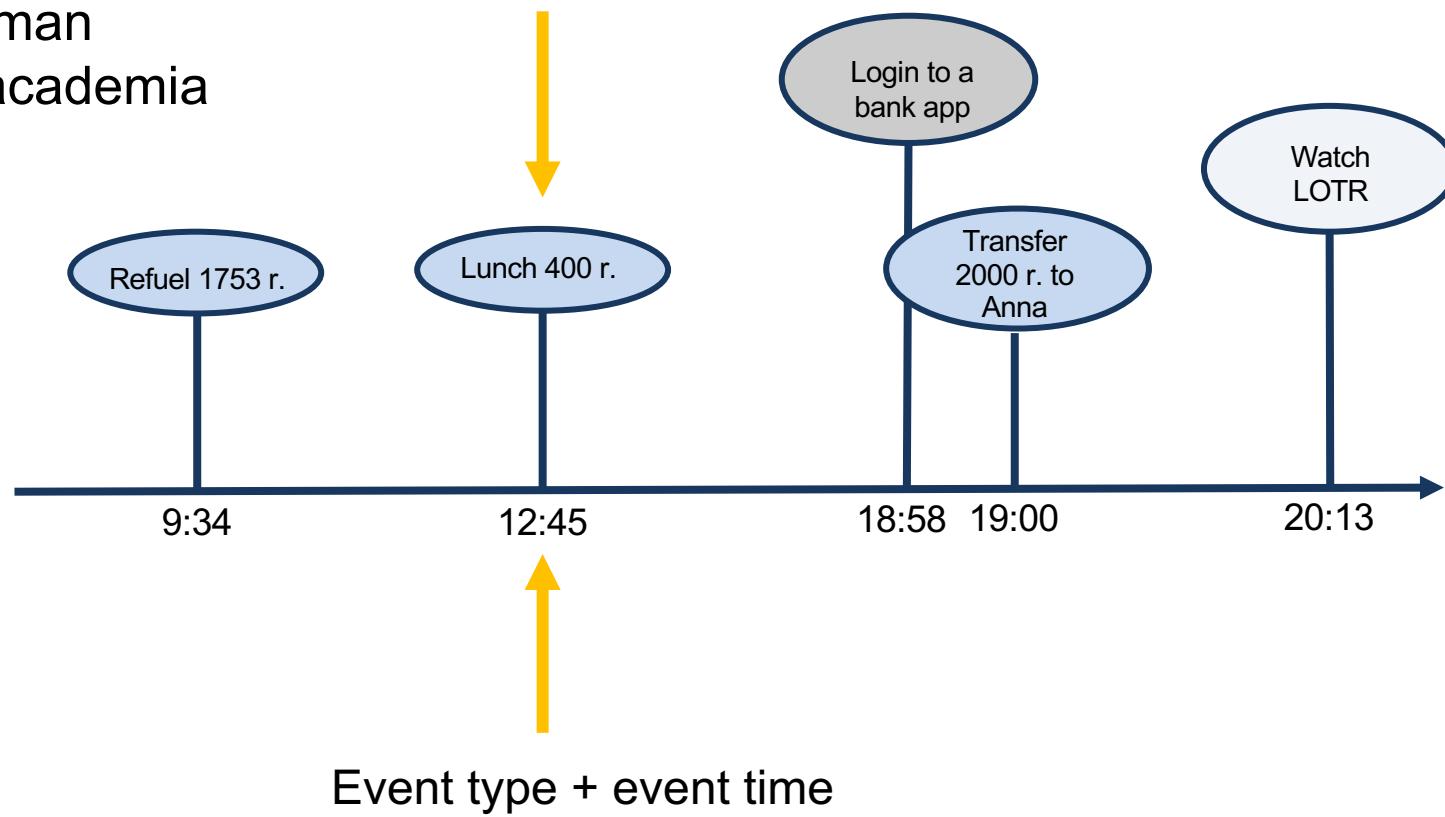
# **Temporal Point Processes (TPPs): applied problems**

# Example: financial transactions are event sequences



Alex, 33, man  
works in academia

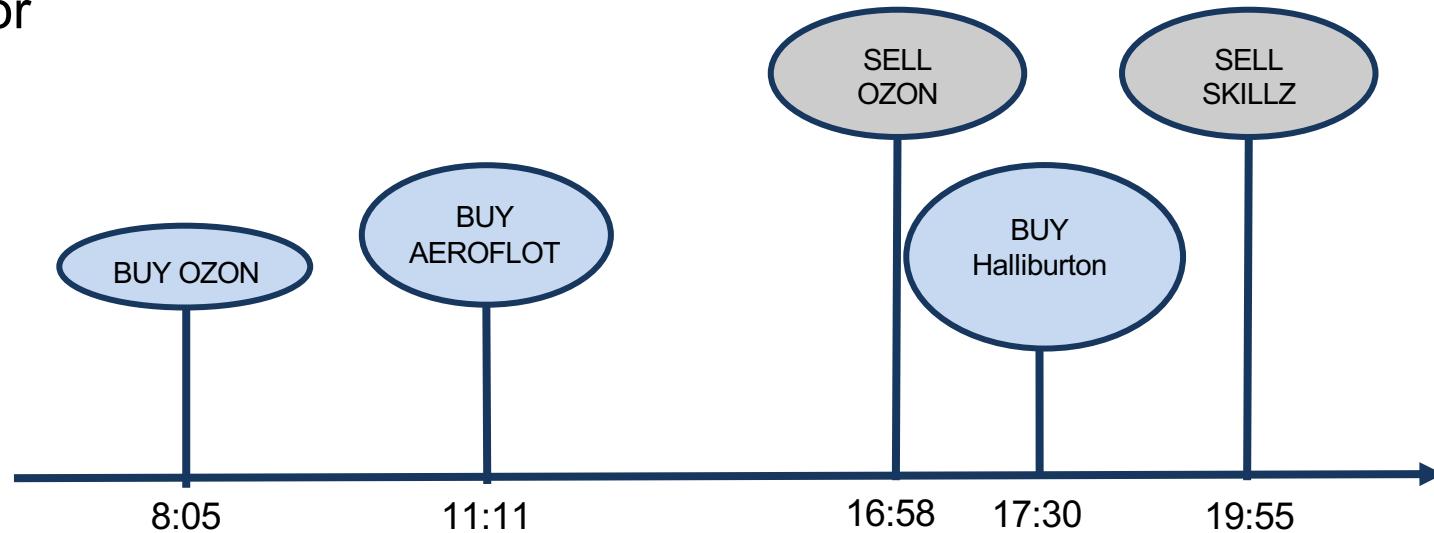
Event features:  
sum, geolocation



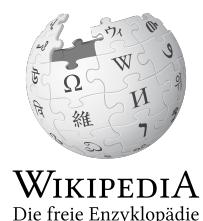
# Example: operations in markets are event sequences



An investor



# More complex example: response history



## Barack Obama

From Wikipedia, the free encyclopedia

*'Barack'* and *'Obama'* redirect here. For his father, see *Barack Obama Sr.* For other uses of *'Barack'*, see *Barack* (disambiguation).

(disambiguation).

Barack Hussein Obama II (born

current President of the United

continental United States. Bo

was president of the Harvard

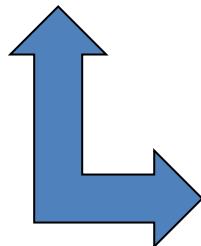
civil rights attorney and taught

representing the 13th District

States House of Representat

## Barack Obama: Revision history

03:41, 28 November 2016 Ranzo (talk | contribs) ... (301,105 bytes) (+18) ... (E  
03:32, 28 November 2016 Xin Deui (talk | contribs) ... (301,087 bytes) (-68) ... (E  
00:57, 28 November 2016 SporkBot (talk | contribs) m ... (301,155 bytes) (-37)  
07:03, 27 November 2016 Saiph121 (talk | contribs) ... (301,192 bytes) (+25) ... (E



03:21, 20 September 2016

is a Kenyan politician

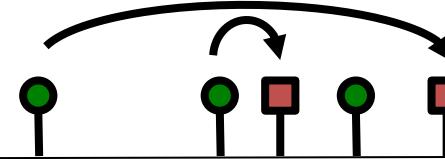
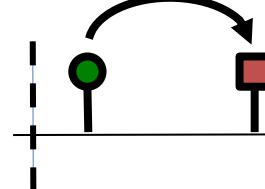


possible vandalism by MLM2016

is an American politician

- Addition
- Refutation

t



Moving to Australia

Working in Australia

Study abroad in Australia

+4

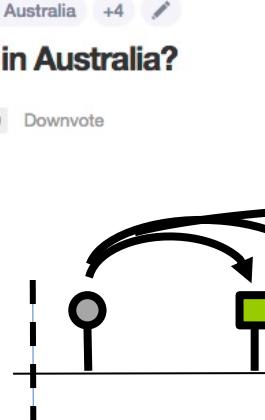
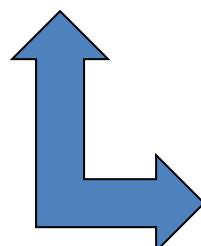
/

Answer

Request

Follow 109 Comment Share 9 Downvote

## What are the pros and cons of living in Australia?



I have studied, worked and lived in Australia as an Intern...

I have experienced this country in all the ways possible, you...  
However, I firmly believe that there are definitely more pros...  
Australia but still I have mentioned below a few challenges and...  
benefits.

Hope it helps! :)

### Possible Challenges

- Language problem for those who don't speak English
- Not having your family and friends around could be challenging as society is more and more connected and thanks to Social Media you can stay in touch a bit easier with them.

Upvote | 150



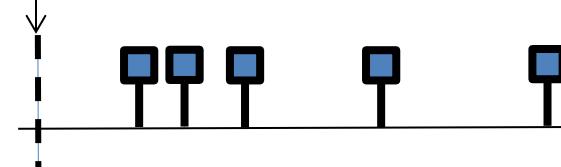
Av M Sharma, Lived in Australia as Migrant, Student, Worker, Business Owner & Family Man

Updated Aug 3

- Question
- Answer

t

↓



Upvote

Skoltech

# Example: development

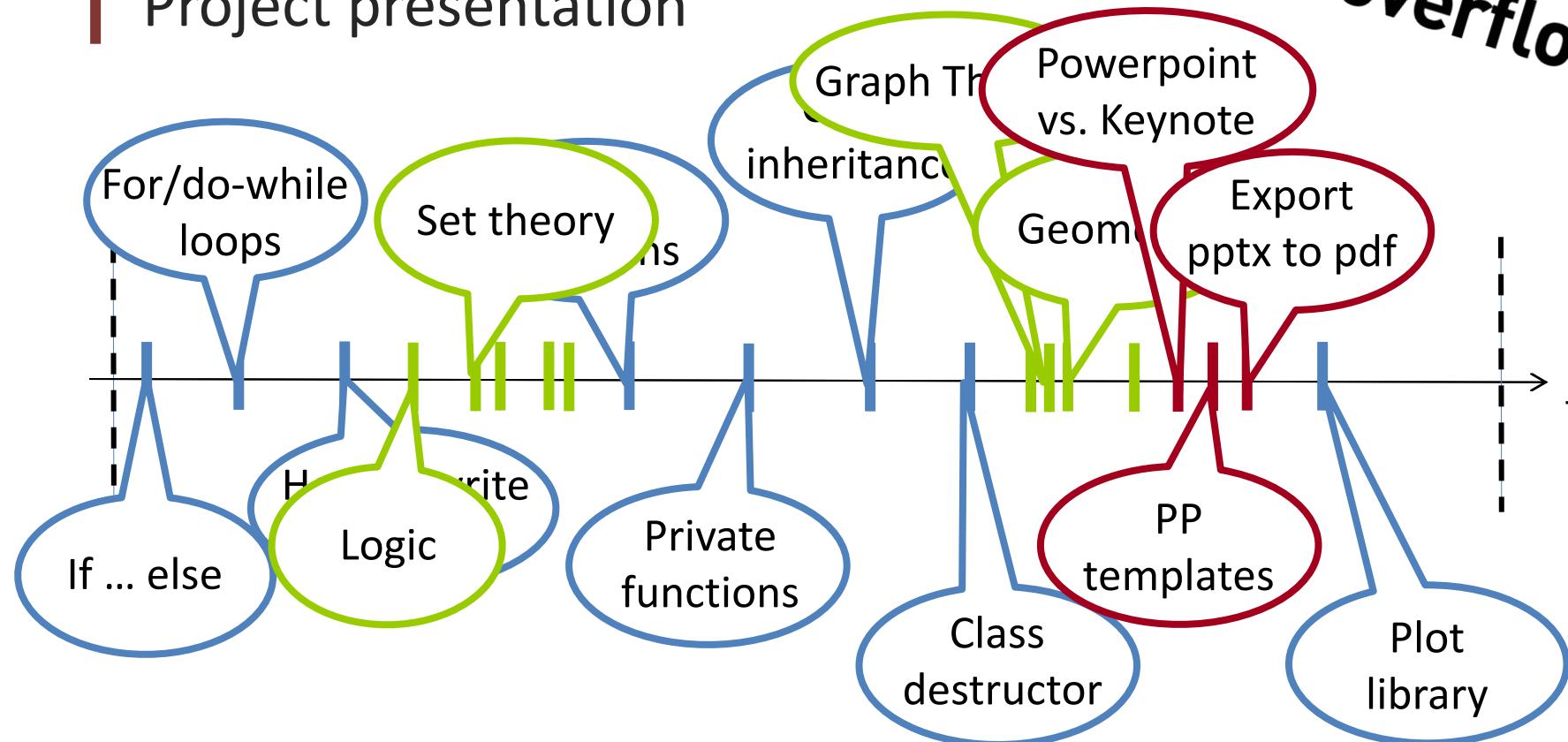


## 1st year computer science student

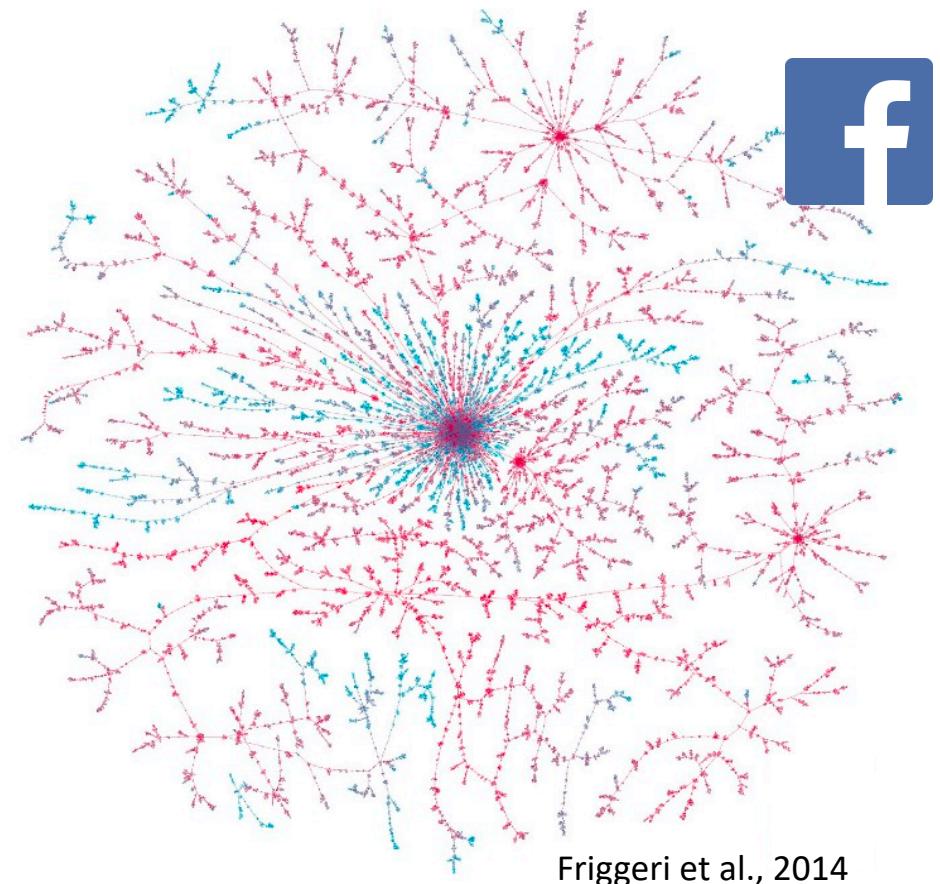
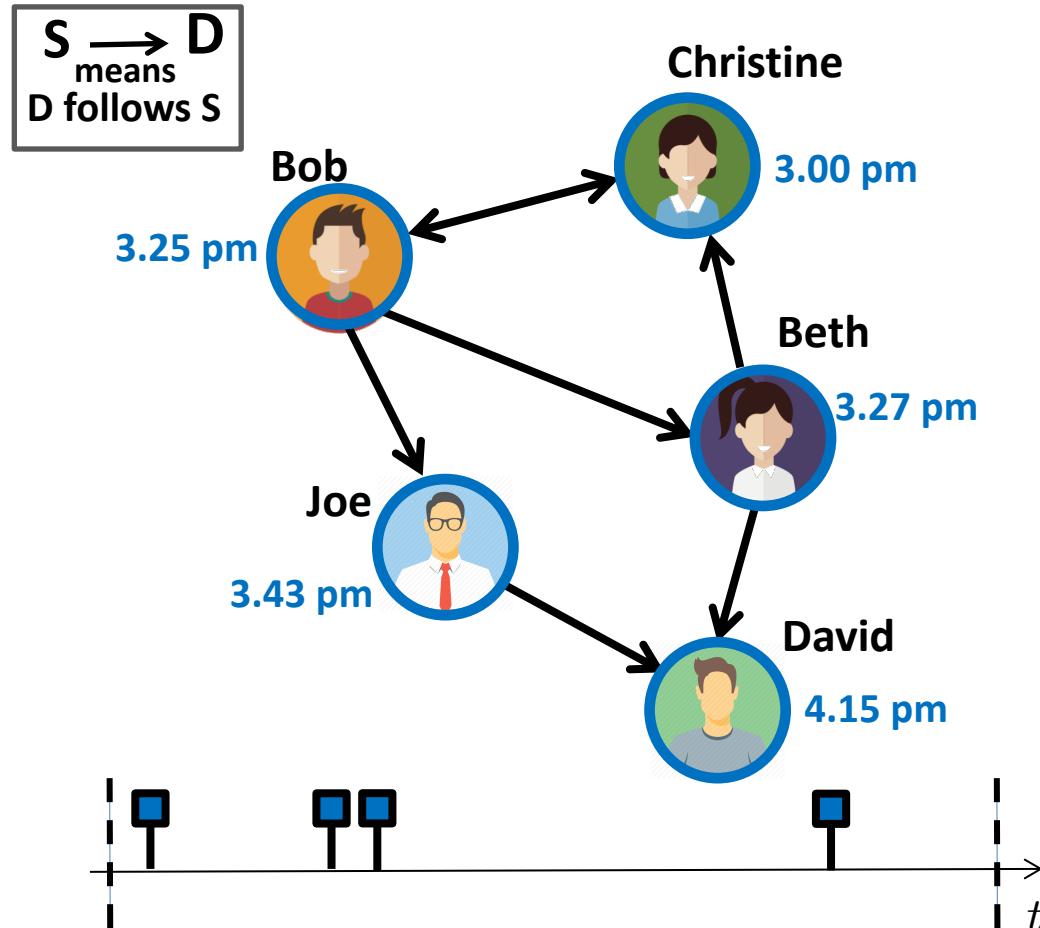
Introduction to programming

Discrete math

Project presentation

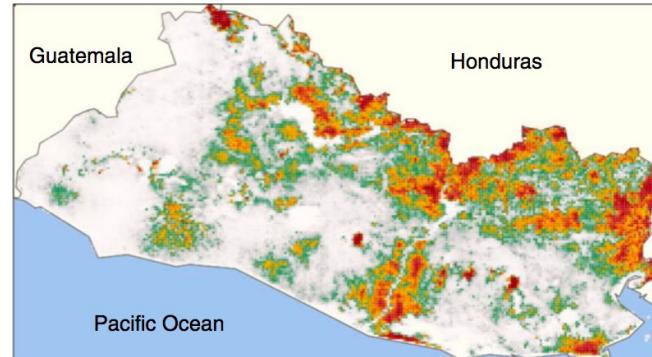
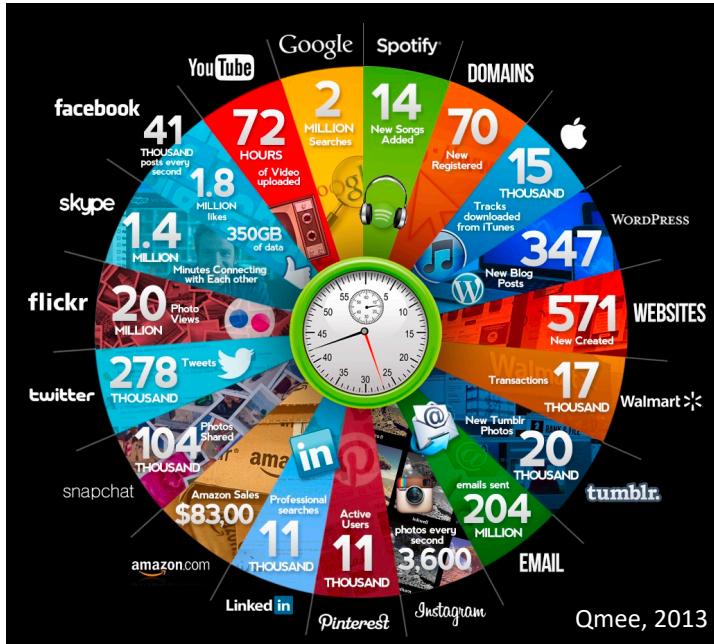


# Example: Information propagation in graphs



These cascades can deliver valuable insights

# Many discrete events in continuous time

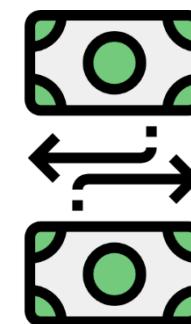
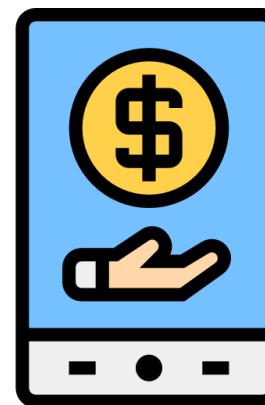


## Disease dynamics

**Clickstreams: Online actions of a user**



**Financial trading**



**Financial transactions**

# Variety of processes behind these events

Events are (noisy) observations of a variety of complex dynamic processes...



Stock trading



Flu spreading



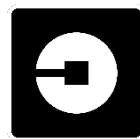
Article creation in Wikipedia



News spread in Twitter



Reviews and sales in Amazon



Ride-sharing requests



A user's reputation in Quora

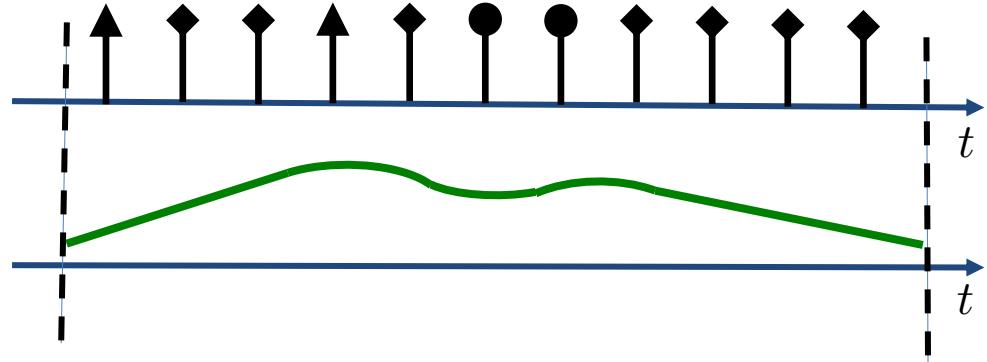
FAST

SLOW

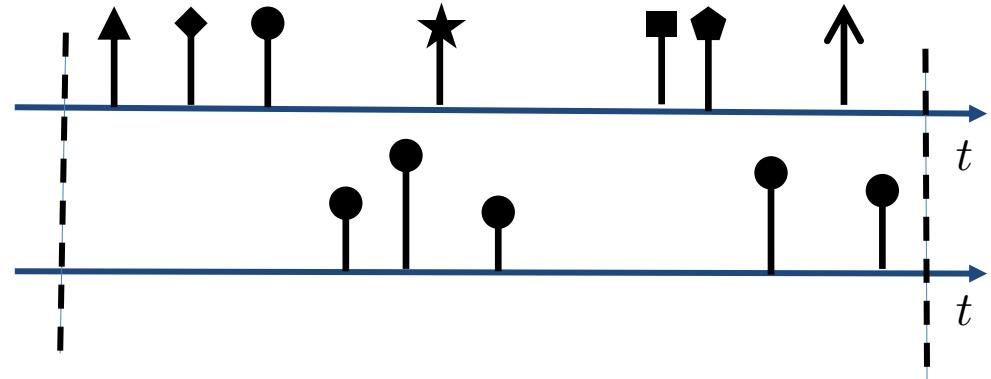


...in a wide range of temporal scales. 9

# Aren't these event traces just time series?

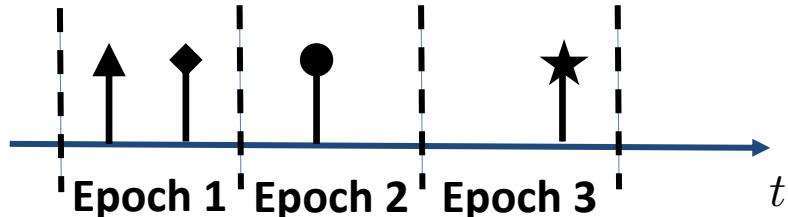


Discrete and continuous times series



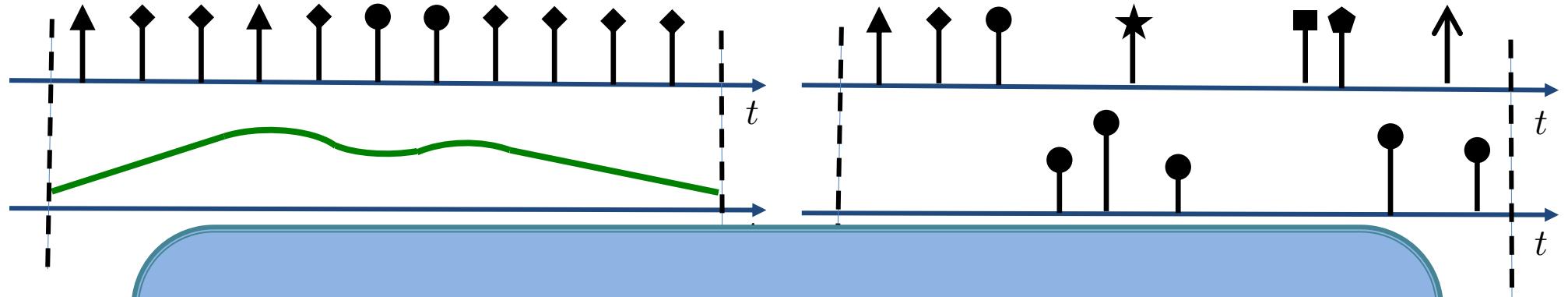
Discrete events in continuous time

What about aggregating events in *epochs*?



- How long is each epoch?
- How to aggregate events per epoch?
- What if no event in one epoch?
- What about time-related queries? 10

# Aren't these event traces just time series?

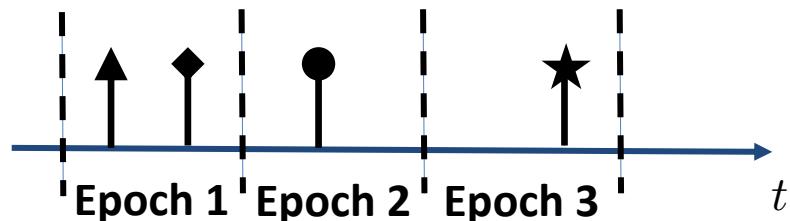


Dis-

The framework of  
temporal point processes  
provides a *native representation*

W

events in epochs?

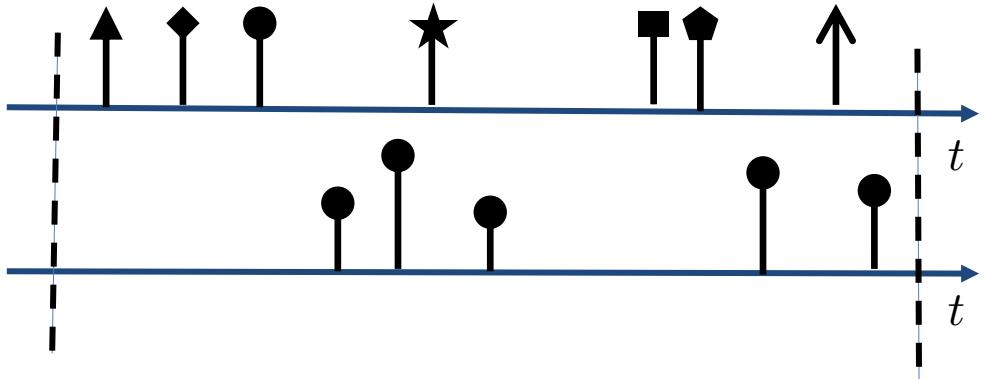


what if no event in one epoch?

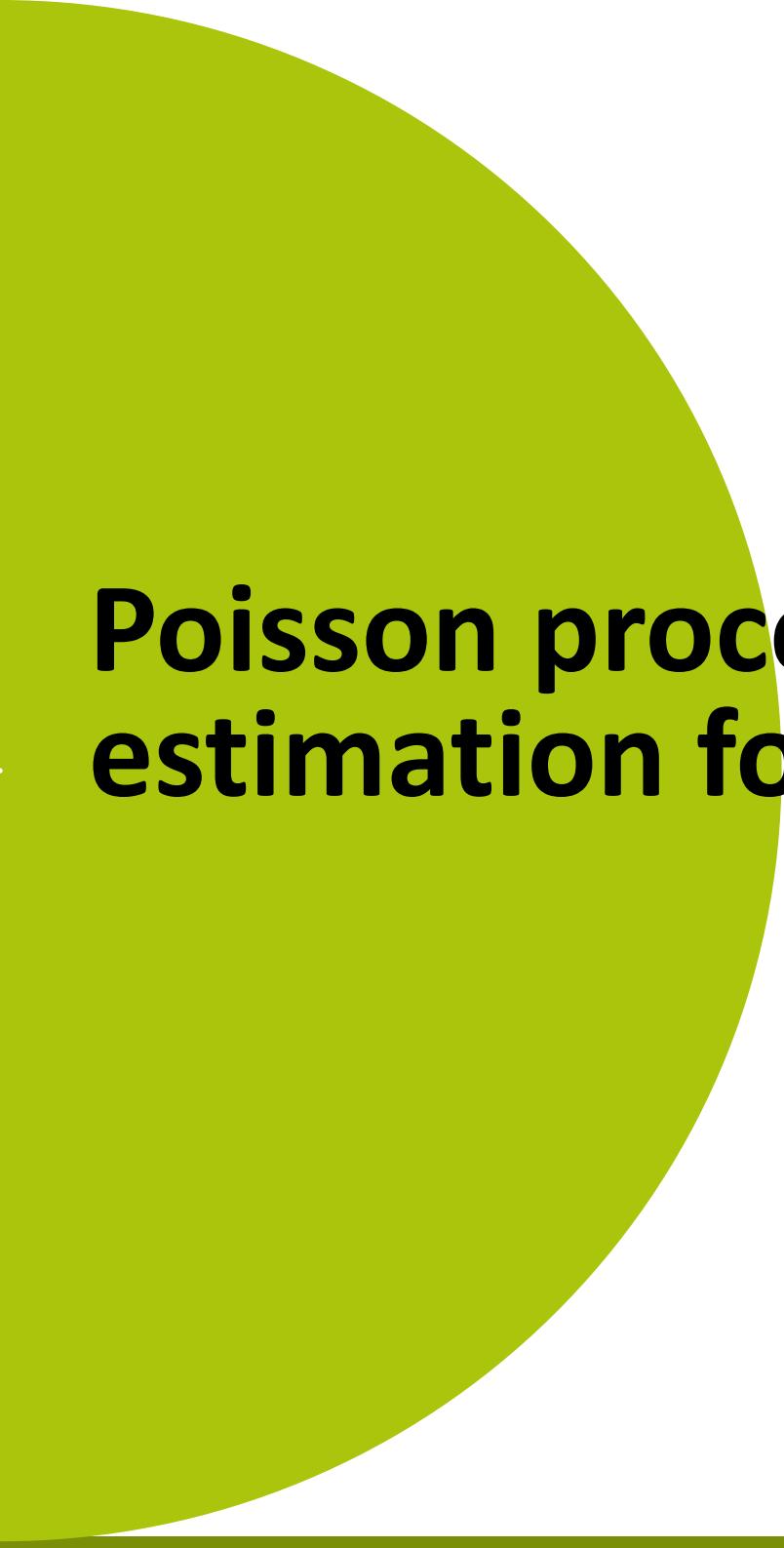
What about time-related queries?

# Problems for event sequences we'll solve!

- Compact description of data: models
- Forecasting/Prediction: distribution for the next event time; event type
- Interpretation: what causes what?
- Control: what should we do?
- Hypothesis testing
- Simulation



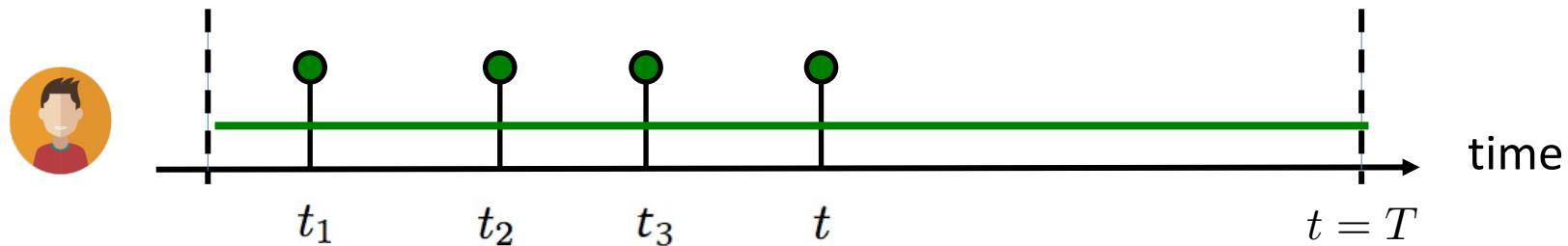
**Point process: generates discrete events in continuous time**



# **Poisson process example and estimation for it; basic terms**

# Poisson process

Coarse approximation of many real-life processes



Intensity of a Poisson process is constant:

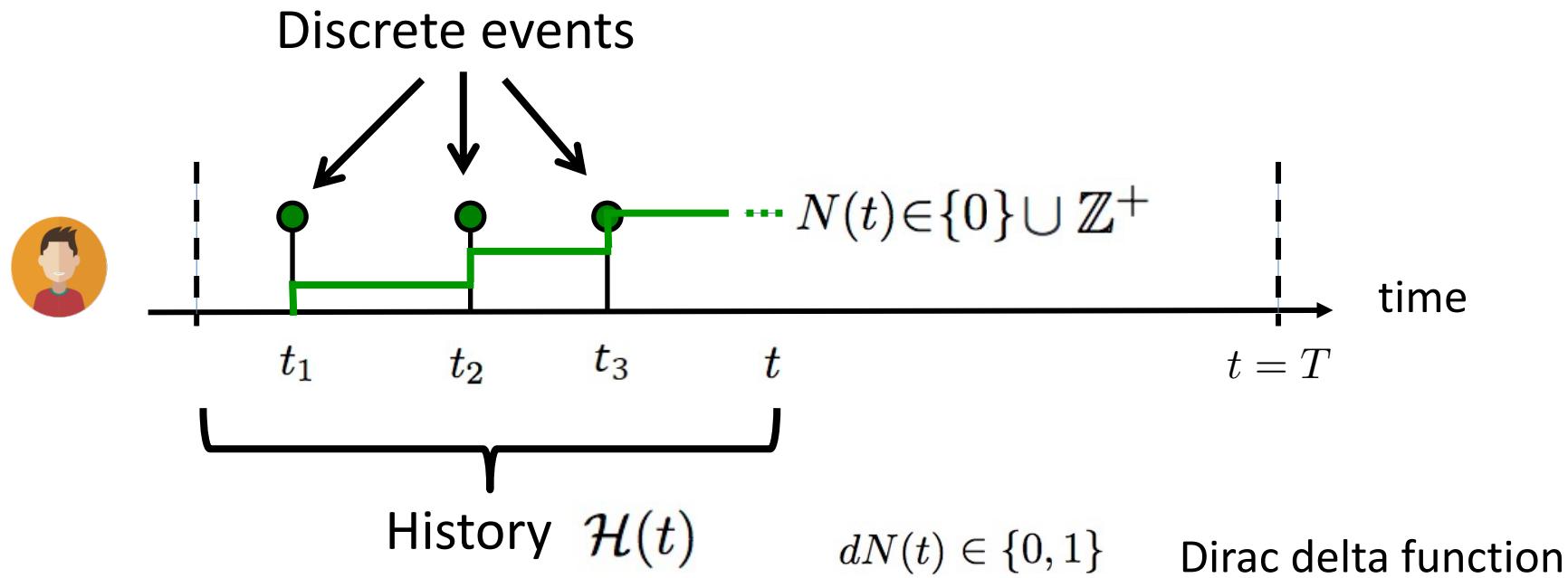
$$\lambda^*(t) = \mu$$

Observations:

1. Intensity independent of history
2. Uniformly random occurrence
3. Time between events follows exponential distribution

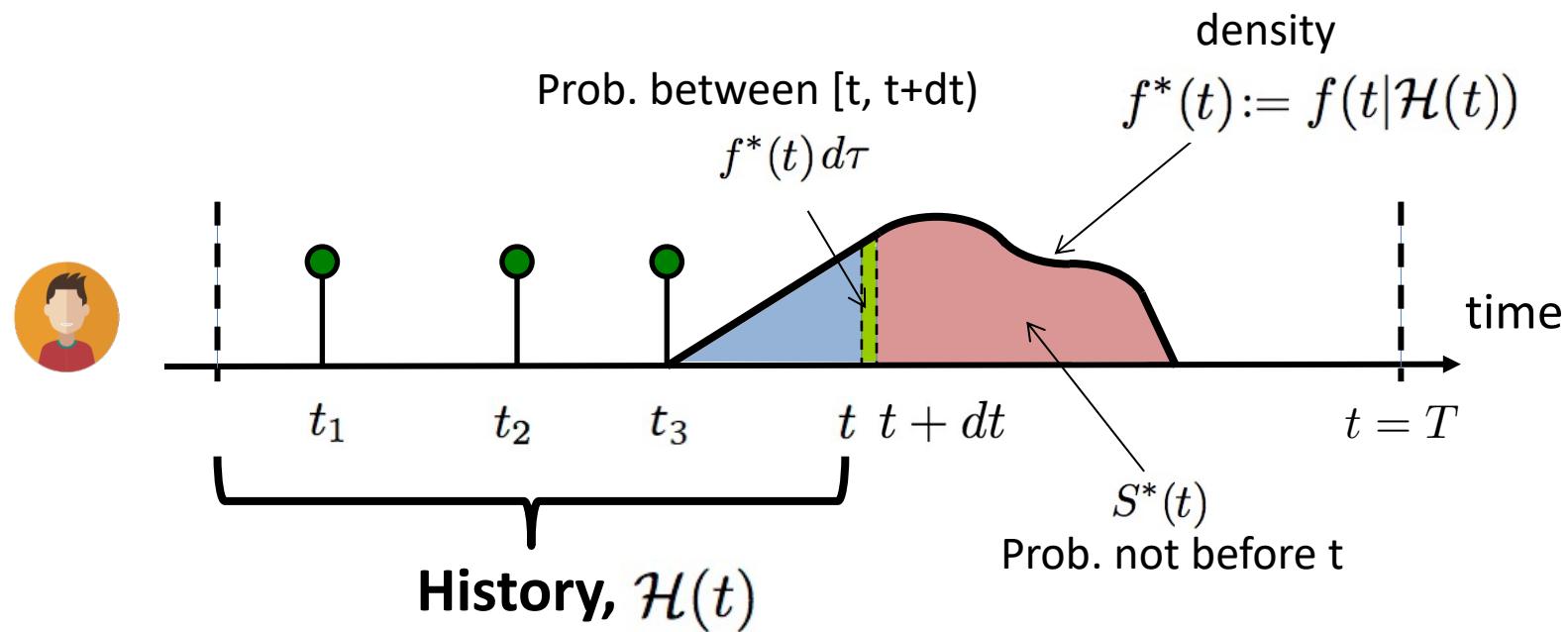
# Defintion: Temporal point processes

A random process whose realization consists of discrete events localized in time  $\mathcal{H} = \{t_i\}$

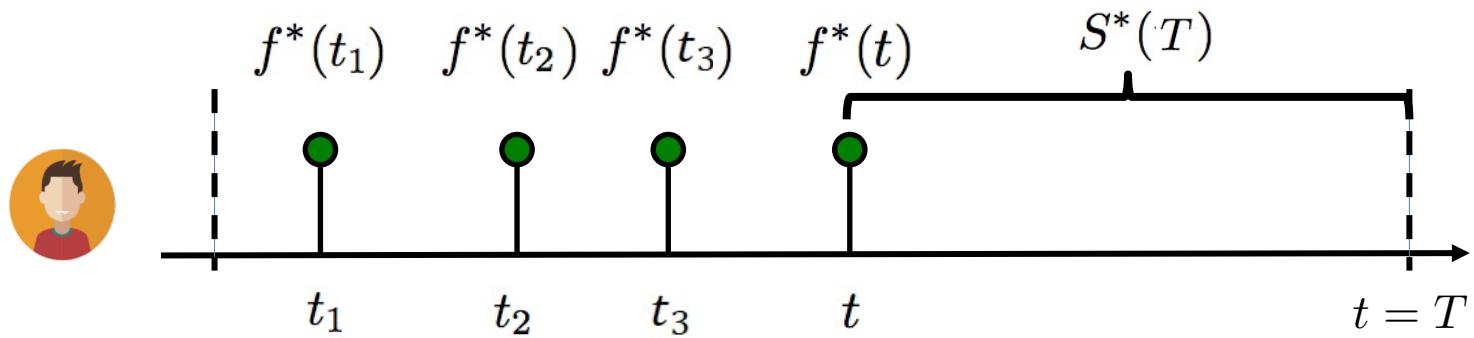


Formally:  $N(t) = \int_0^t dN(s) \Rightarrow dN(t) = \sum_{t_i \in \mathcal{H}} \delta(t - t_i) dt$

# Distribution we are looking for the next event time

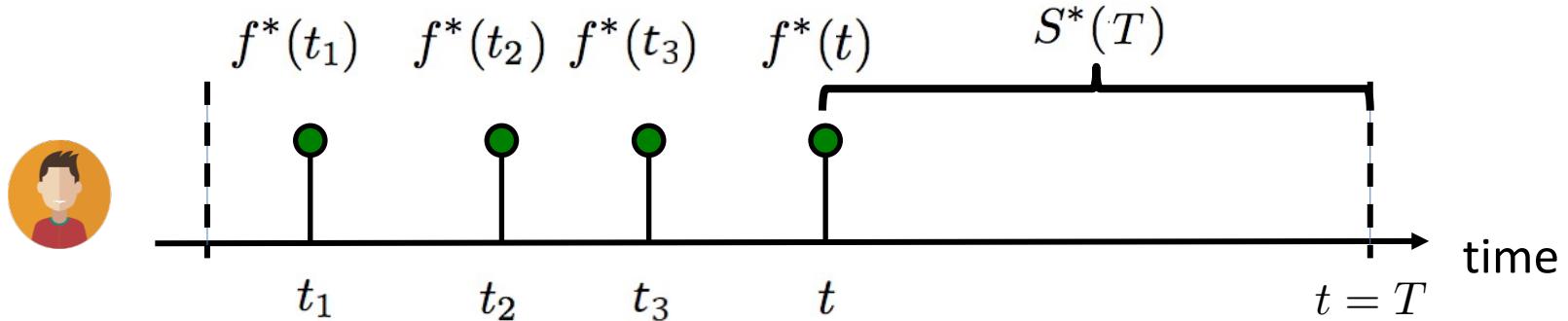


# We can write and optimize a likelihood function



Likelihood of a series:  $f^*(t_1) \ f^*(t_2) \ f^*(t_3) \ f^*(t) \ S^*(T)$

# Density parametrization is hard

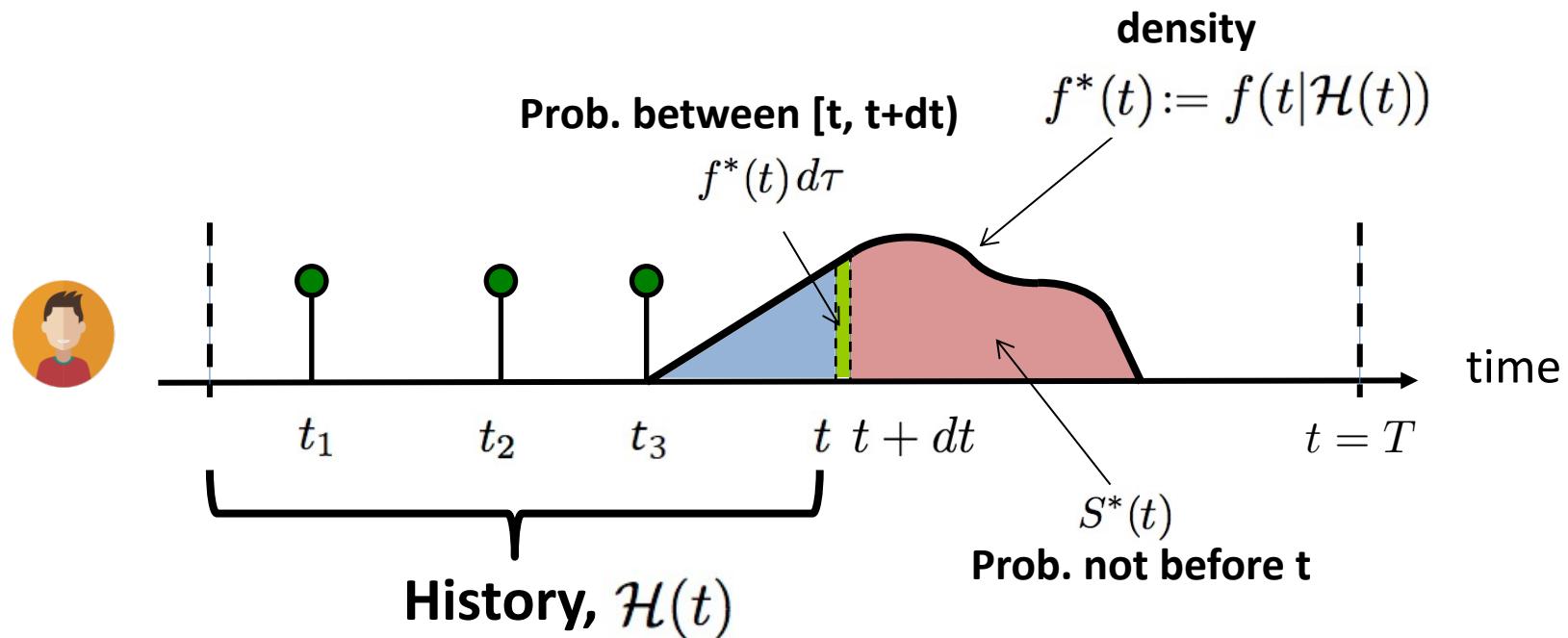


$$\begin{array}{c} f^*(t_1) \quad f^*(t_2) \quad f^*(t_3) \quad f^*(t) \quad S^*(T) \\ \uparrow \qquad \uparrow \qquad \uparrow \qquad \uparrow \qquad \swarrow \\ \frac{\exp\langle w, \psi^*(t_1) \rangle}{Z} \quad \frac{\exp\langle w, \psi^*(t_2) \rangle}{Z} \quad \frac{\exp\langle w, \psi^*(t_3) \rangle}{Z} \quad \frac{\exp\langle w, \psi^*(t) \rangle}{Z} \quad 1 - \int_t^T \frac{\exp\langle w, \psi^*(\tau) \rangle}{Z} d\tau \end{array}$$

It is difficult for model design and interpretability:

1. Densities need to integrate to 1 (i.e., partition function)
2. Difficult to combine timelines

# Intensity function is an alternative



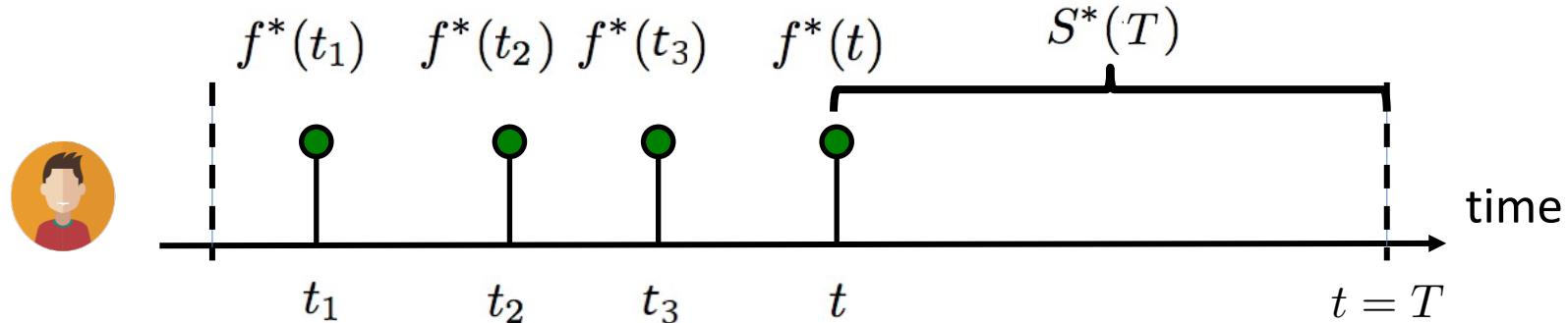
**Intensity:**

Probability between  $[t, t+dt]$  but not before  $t$

$$\lambda^*(t)dt = \frac{f^*(t)dt}{S^*(t)} \geq 0 \quad \rightarrow \quad \lambda^*(t)dt = \mathbb{E}[dN(t)|\mathcal{H}(t)]$$

**Note:**  $\lambda^*(t)$  is a rate = # of events / unit of time

# Likelihood for intensity

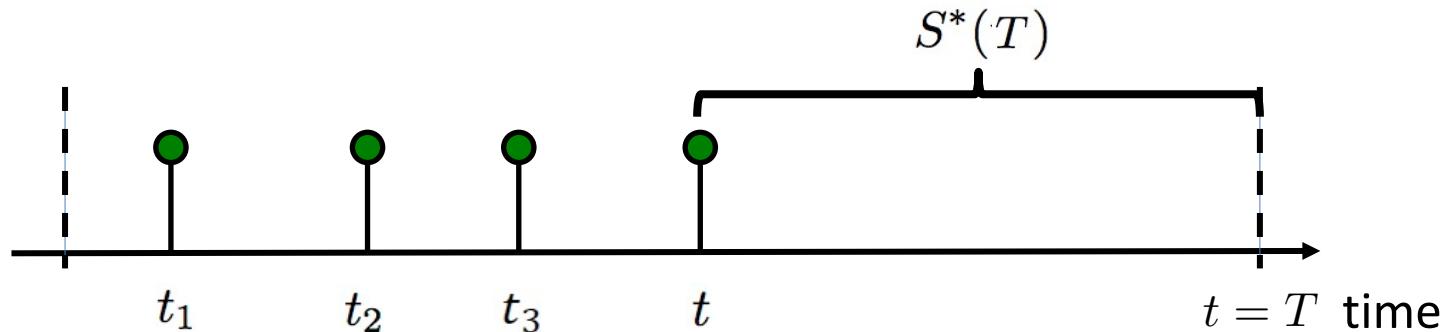


$$\lambda^*(t_1) \lambda^*(t_2) \lambda^*(t_3) \lambda^*(t) \exp\left(-\int_0^T \lambda^*(\tau) d\tau\right)$$
$$\langle w, \phi^*(t_1) \rangle \quad \langle w, \phi^*(t_2) \rangle \quad \langle w, \phi^*(t_3) \rangle \quad \langle w, \phi^*(t) \rangle$$
$$\exp\left(-\int_0^T \langle w, \phi^*(\tau) \rangle d\tau\right)$$

Suitable for model design and interpretable:

1. Intensities only need to be nonnegative
2. Easy to combine timelines

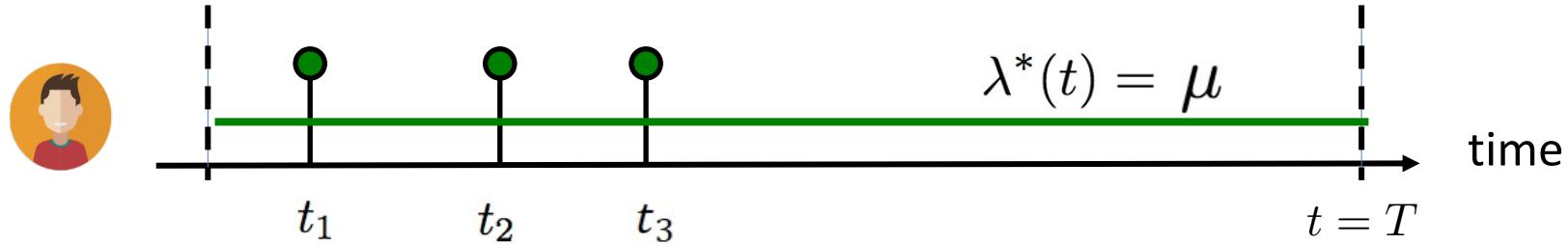
# Log likelihood via intensity and density



$$L = f(t_1|\mathcal{H}_0)f(t_2|\mathcal{H}_{t_1}) \cdots f(t_n|\mathcal{H}_{t_{n-1}})(1 - F(T|\mathcal{H}_{t_n}))$$

$$\begin{aligned} L &= \left( \prod_{i=1}^n f(t_i|\mathcal{H}_{t_{i-1}}) \right) \frac{f(T|\mathcal{H}_{t_n})}{\lambda^*(T)} \\ &= \left( \prod_{i=1}^n \lambda^*(t_i) \exp \left( - \int_{t_{i-1}}^{t_i} \lambda^*(s) ds \right) \right) \exp \left( - \int_{t_n}^T \lambda^*(s) ds \right) \\ &= \left( \prod_{i=1}^n \lambda^*(t_i) \right) \exp \left( - \int_0^T \lambda^*(s) ds \right), \end{aligned}$$

# Fitting & sampling for a Poisson process



Fitting by maximum likelihood:

$$\mu^* = \underset{\mu}{\operatorname{argmax}} \ 3 \log \mu - \mu T = \frac{3}{T}$$

Sampling using inversion sampling:

$$t \sim \mu \exp(-\mu(t - t_3))$$

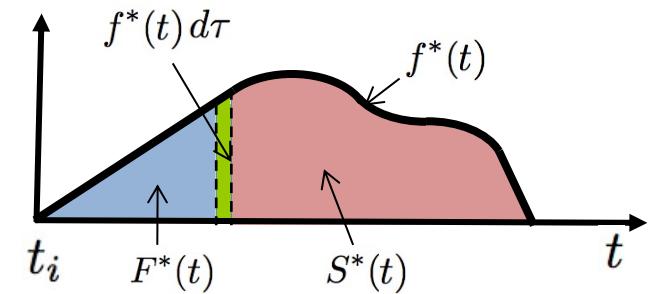
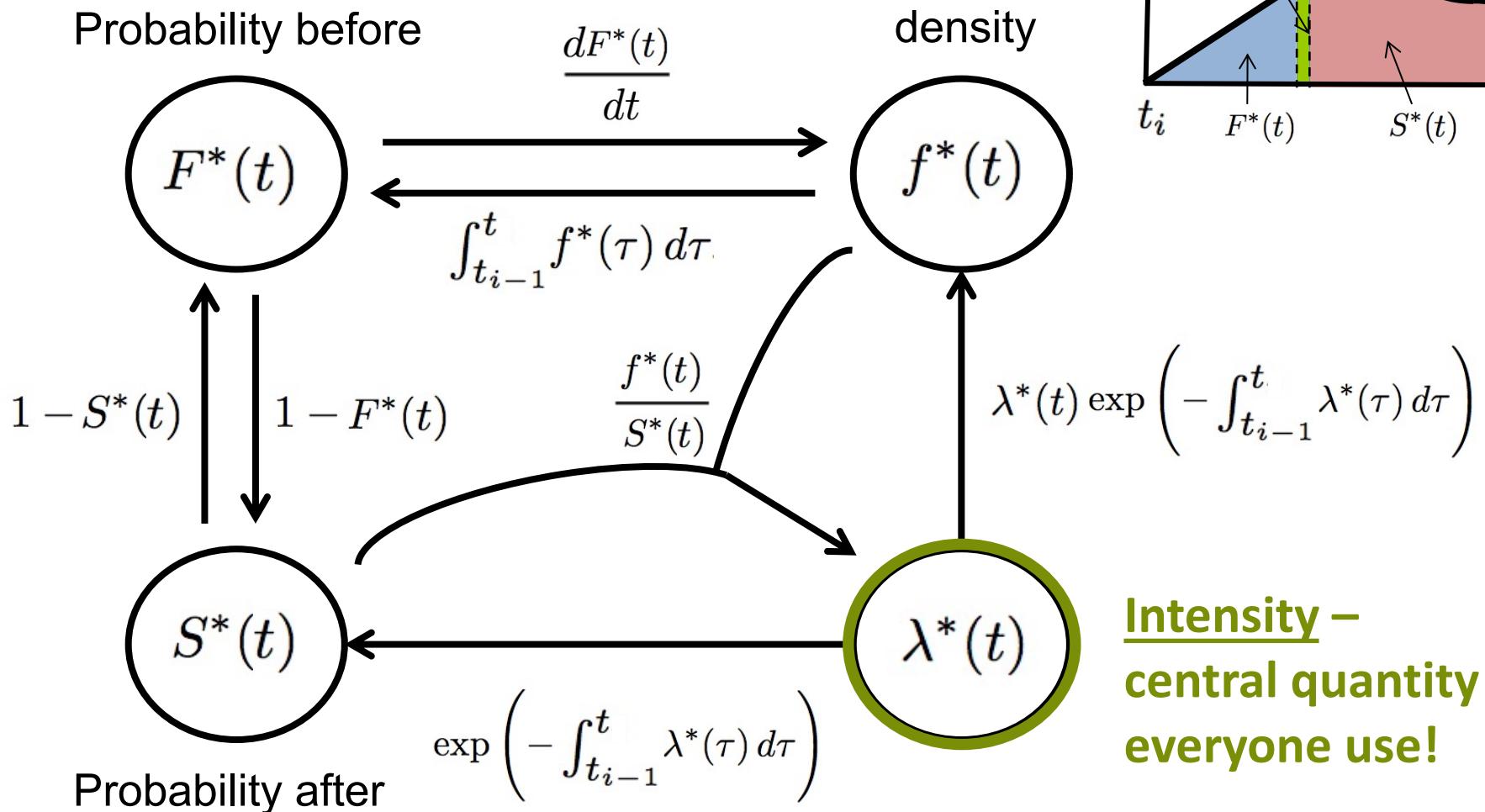


$$t = -\frac{1}{\mu} \log(1 - u) + t_3$$

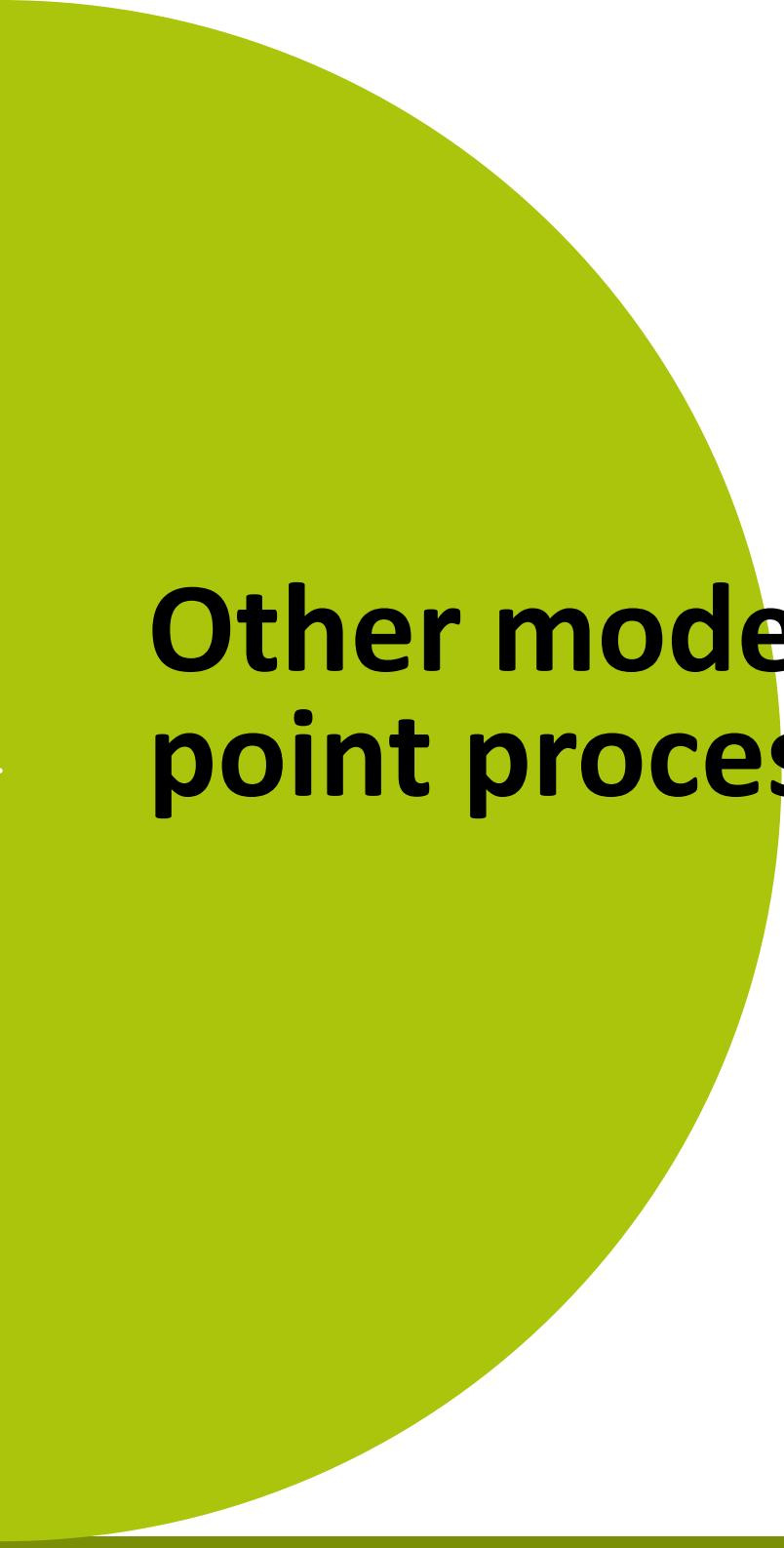
$$f_t^*(t)$$

$$F_t^{-1}(u)$$

# Relation between $f^*$ , $F^*$ , $S^*$ , $\lambda^*$



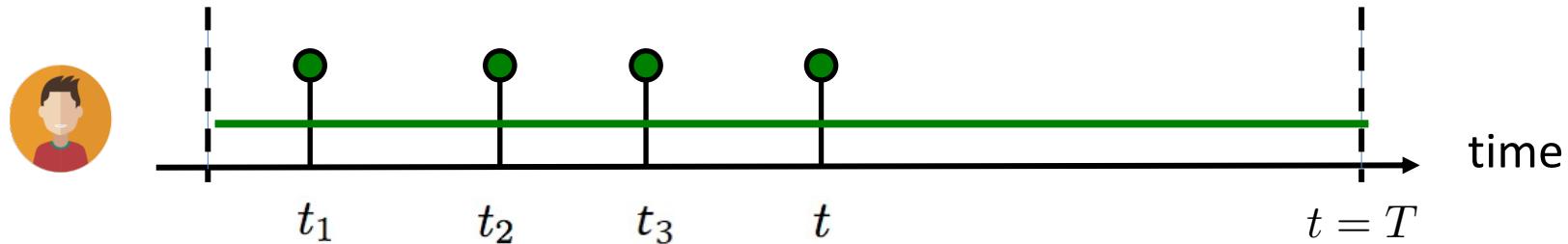
**Intensity –  
central quantity  
everyone use!**



# **Other models for temporal point processes**

# Poisson process

Coarse approximation of many real-life processes



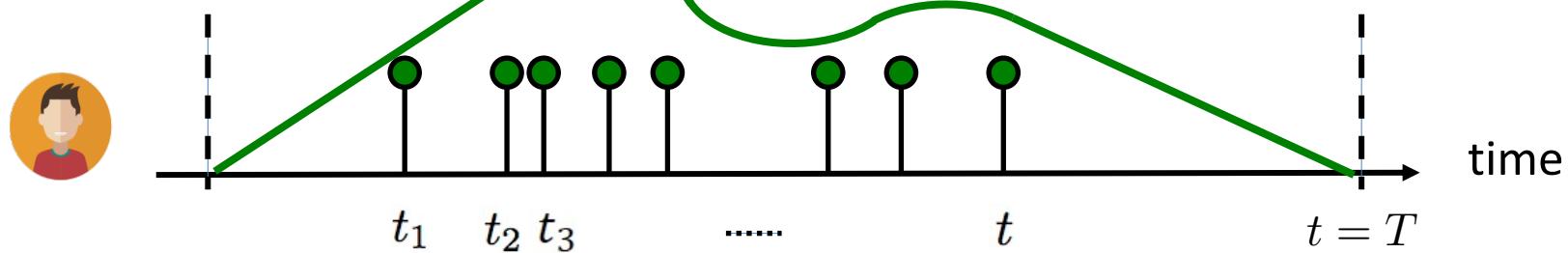
Intensity of a Poisson process

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

Observations:

1. Intensity independent of history
2. Uniformly random occurrence
3. Time interval follows exponential distribution

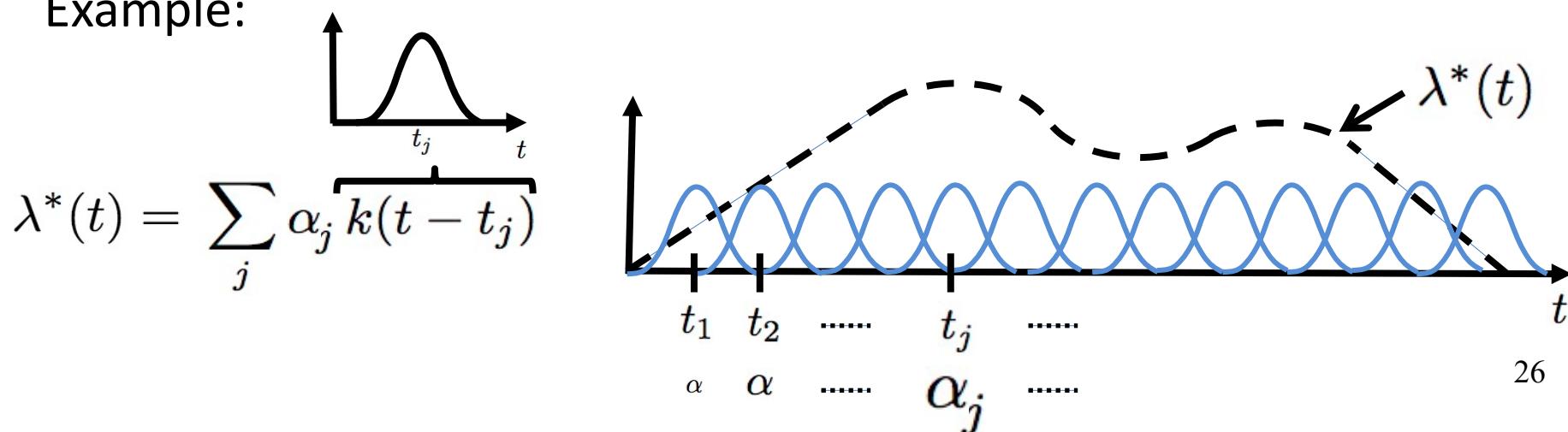
# Inhomogeneous Poisson process



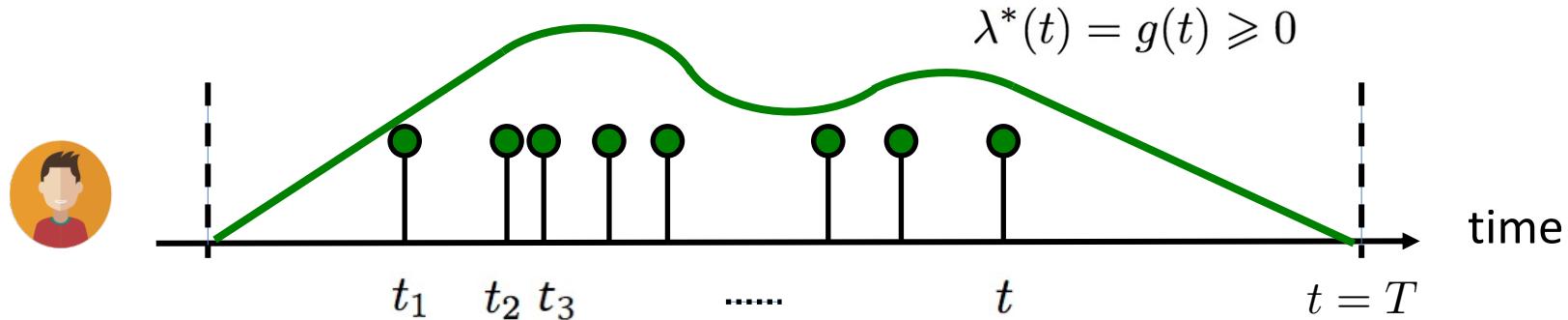
Intensity of an inhomogeneous Poisson process

$$\lambda^*(t) = g(t) \geq 0 \quad -\text{Independent of history}$$

Example:



# Fitting from inhomogeneous Poisson



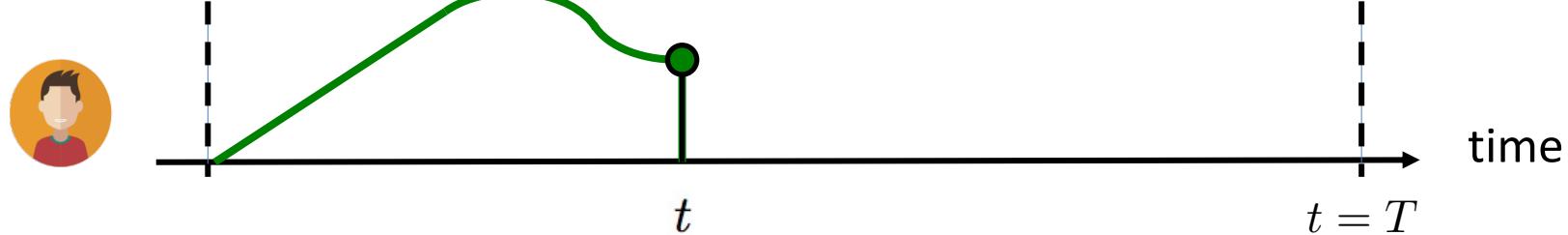
Fitting by maximum log-likelihood:

$$\underset{g(t)}{\text{maximize}} \quad \sum_{i=1}^n \log g(t_i) - \int_0^T g(\tau) d\tau$$

Idea: we have additional features, so we can use a generalized linear model for it

Intensity is  $g(t) = g(\mathbf{x}_t) = \exp(\mathbf{x}_t^T \mathbf{w})$

# Terminating (or survival) process



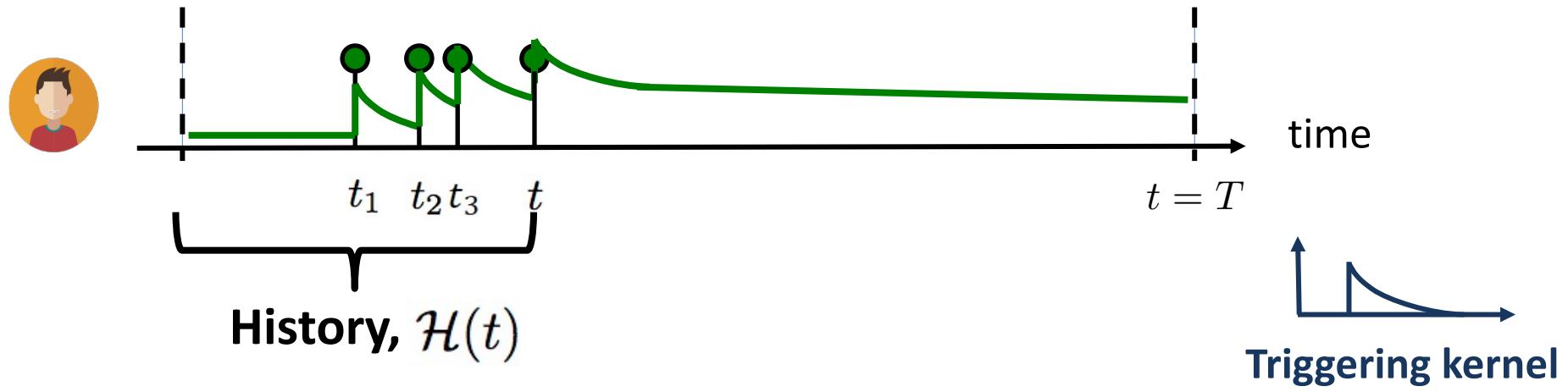
Intensity of a terminating (or survival) process

$$\lambda^*(t) = g^*(t)(1 - N(t)) \geq 0$$

Observations:

1. Limited number of occurrences
2. Hazard function in actuarial science

# Self-exciting Hawkes process



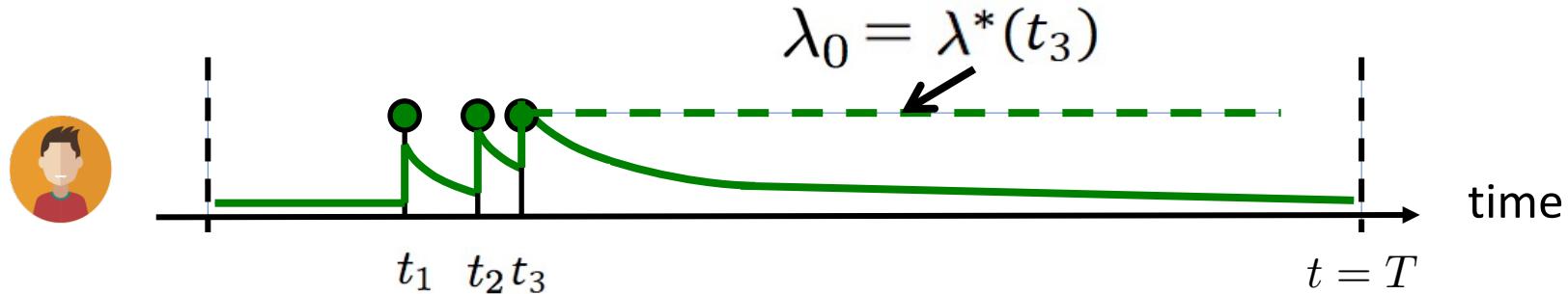
Intensity of self-exciting  
(or Hawkes) process:

$$\begin{aligned}\lambda^*(t) &= \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i) \\ &= \mu + \alpha \kappa_\omega(t) \star dN(t)\end{aligned}$$

Observations:

1. Clustered (or bursty) occurrence of events
2. Intensity is stochastic and history dependent

# Fitting a Hawkes process from a recorded timeline



Fitting by maximum likelihood:

$$\underset{\mu, \alpha}{\text{maximize}} \quad \sum_{i=1}^n \log \lambda^*(t_i) - \int_0^T \lambda^*(\tau) d\tau \quad \left. \right\} \begin{array}{l} \text{The max. likelihood} \\ \text{is jointly convex} \\ \text{in } \mu \text{ and } \alpha \end{array}$$

# Summary

**Building blocks to represent different dynamic processes:**

Poisson processes:

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

Inhomogeneous Poisson processes:

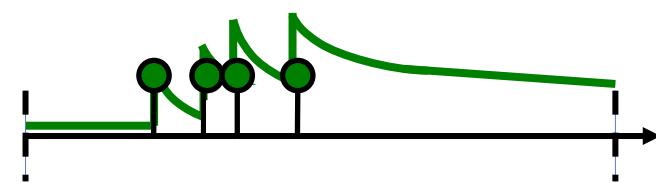
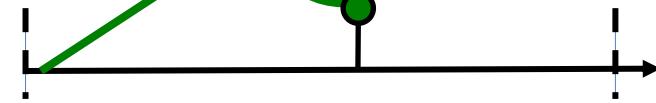
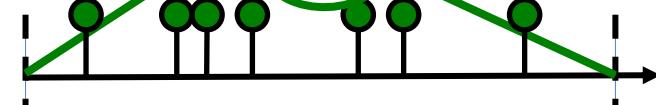
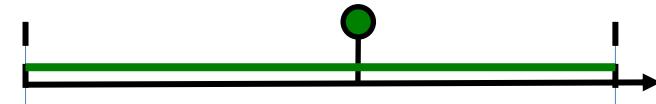
$$\lambda^*(t) = g(t)$$

Terminating point processes:

$$\lambda^*(t) = g^*(t)(1 - N(t))$$

Self-exciting point processes:

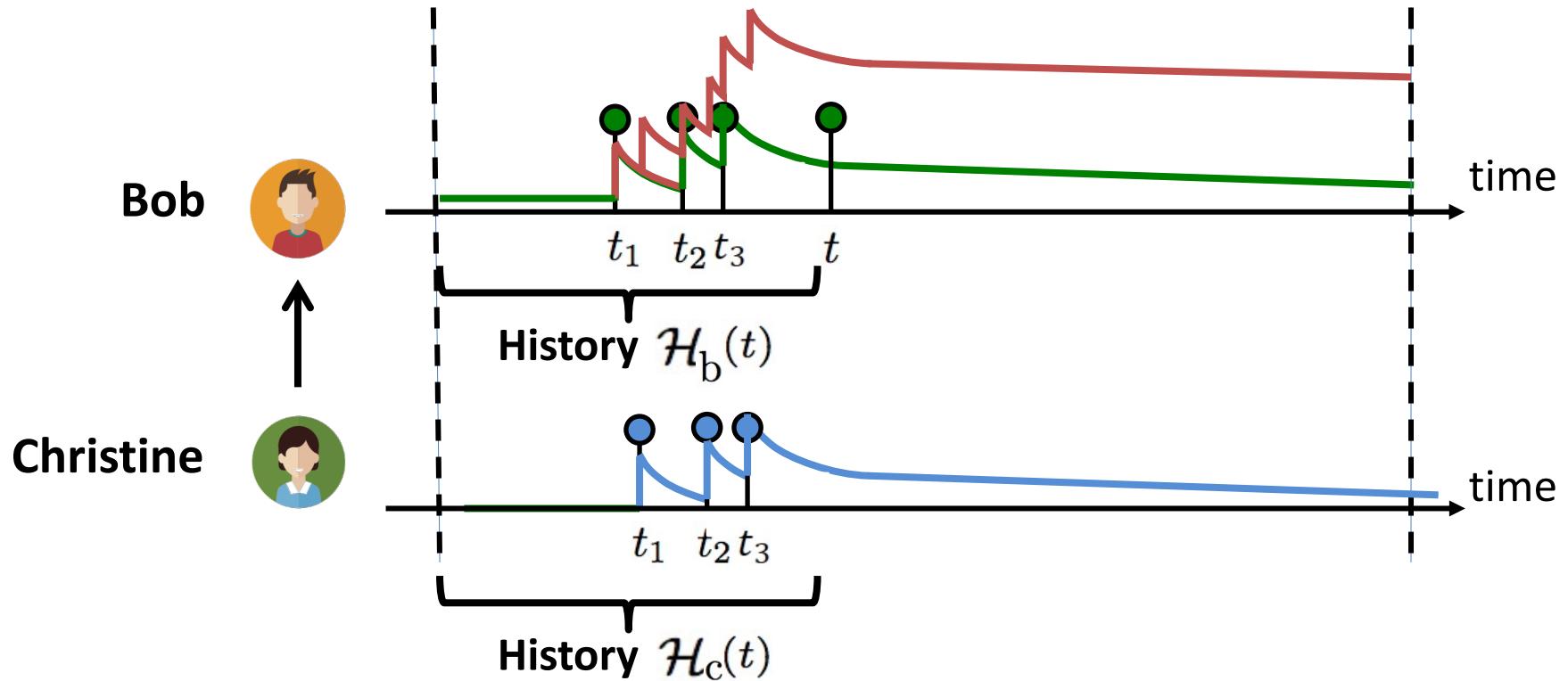
$$\lambda^*(t) = \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i)$$





# **Mutually excited point processes**

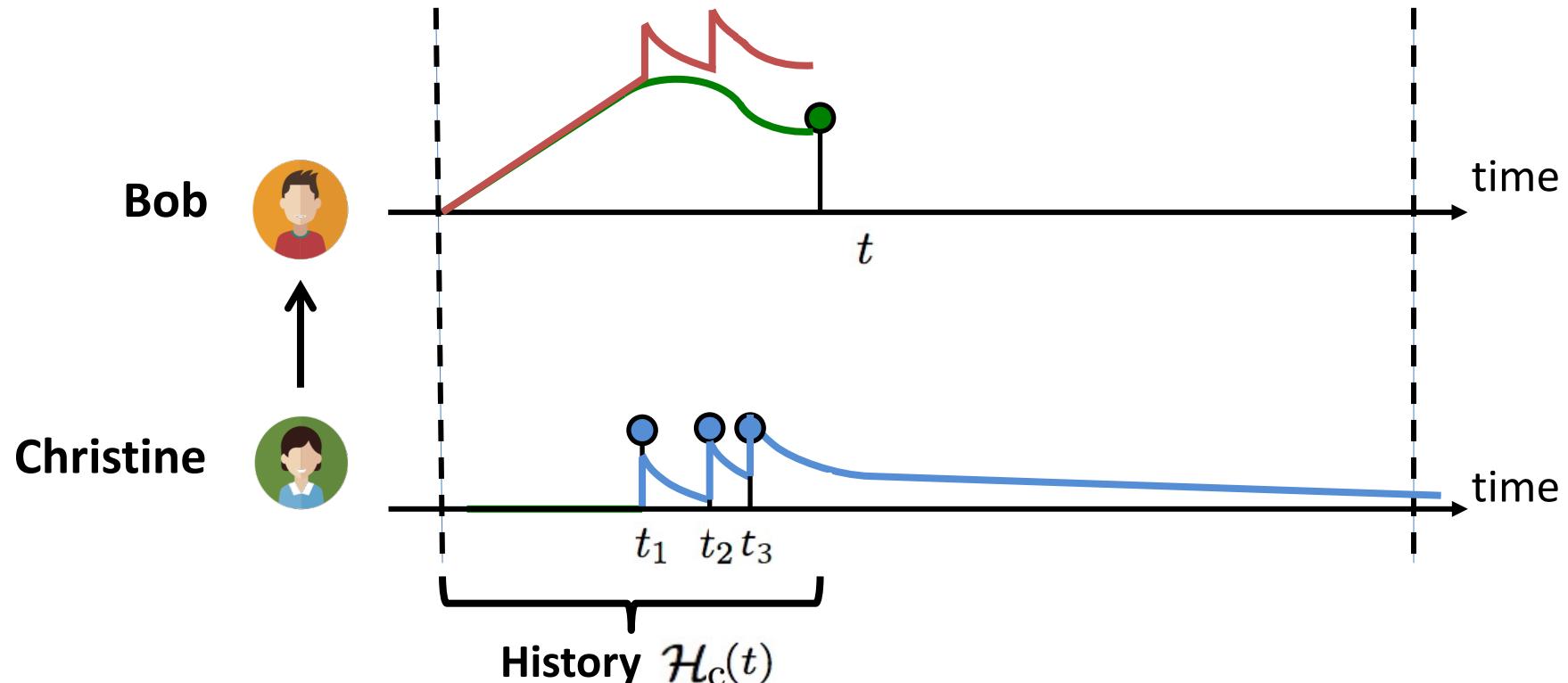
# Mutually exciting process



Clustered occurrence affected by neighbors

$$\begin{aligned}\lambda^*(t) = & \mu + \alpha \sum_{t_i \in \mathcal{H}_b(t)} \kappa_\omega(t - t_i) \\ & + \beta \sum_{t_i \in \mathcal{H}_c(t)} \kappa_\omega(t - t_i)\end{aligned}$$

# Mutually exciting terminating process



Clustered occurrence affected by neighbors

$$\lambda^*(t) = (1 - N(t)) \left( g(t) + \beta \sum_{t_i \in \mathcal{H}_c(t)} \kappa_\omega(t - t_i) \right)$$

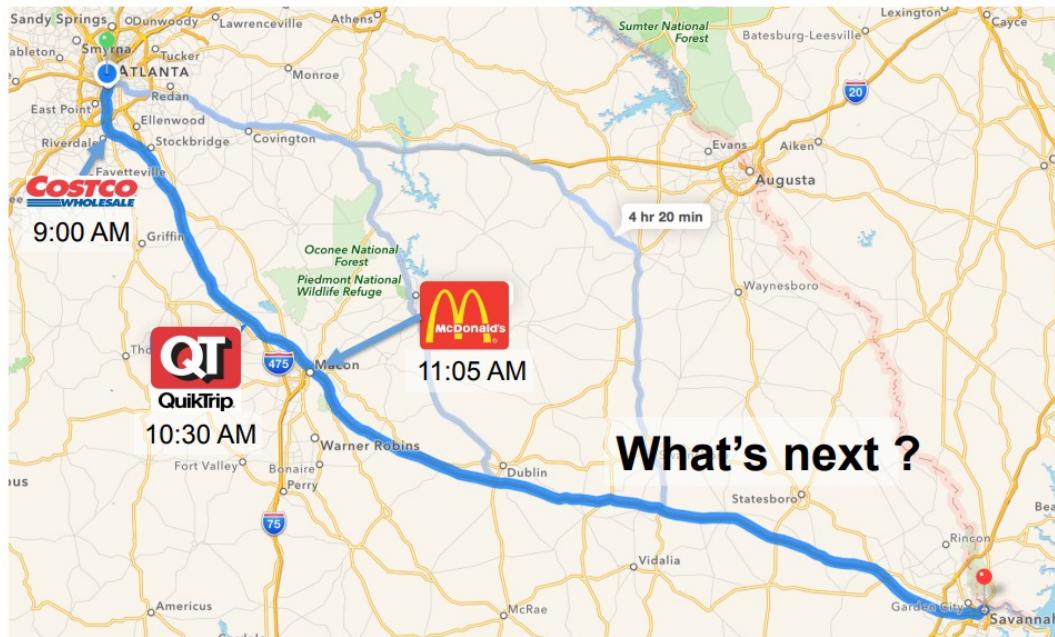


# **Marked temporal point processes**

# Marked temporal point processes

Marked temporal point process:

A random process whose realization consists of discrete *marked* events localized in time

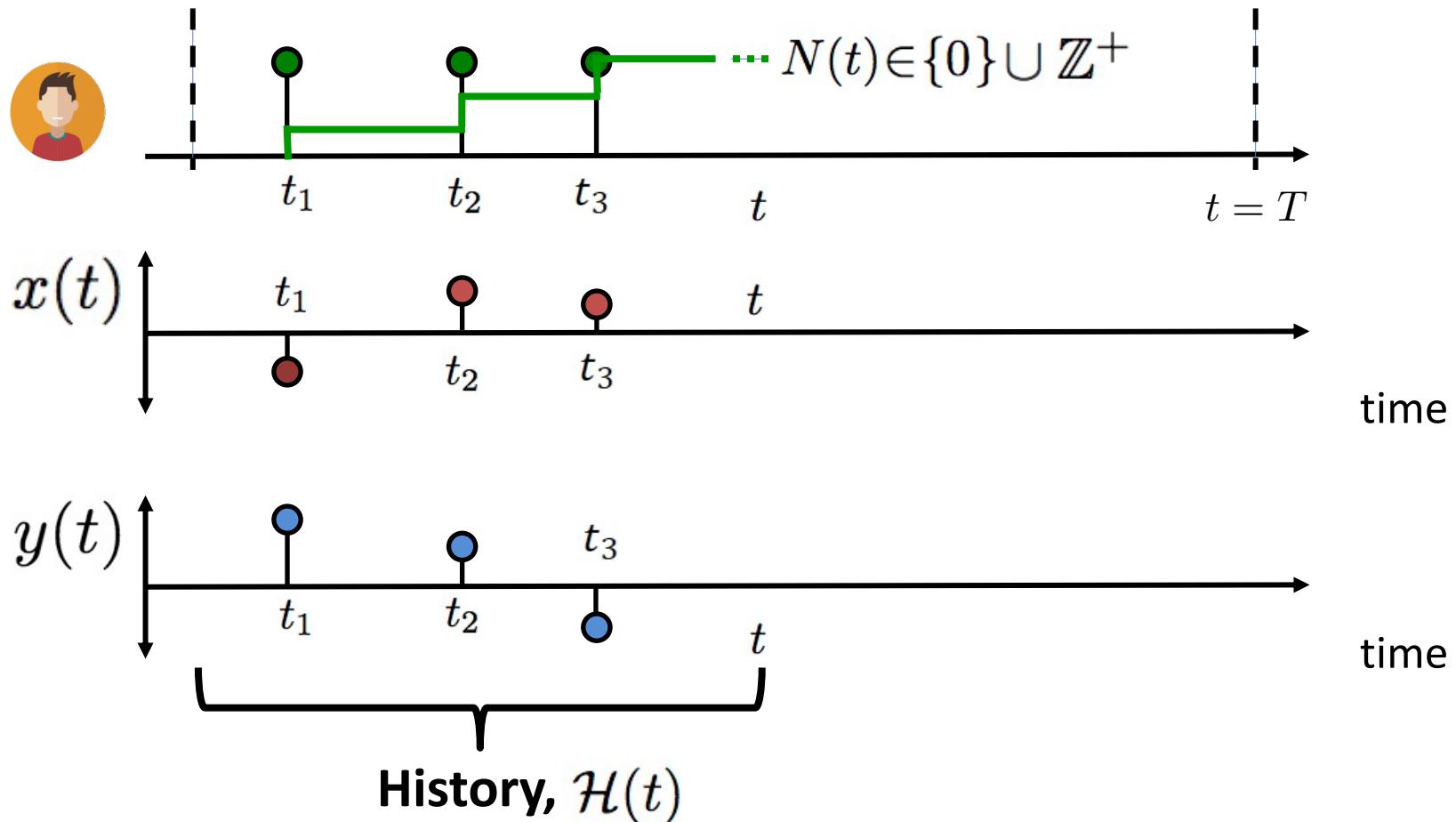


Given the trace of past locations and time, can we predict the location and time of the next stop?

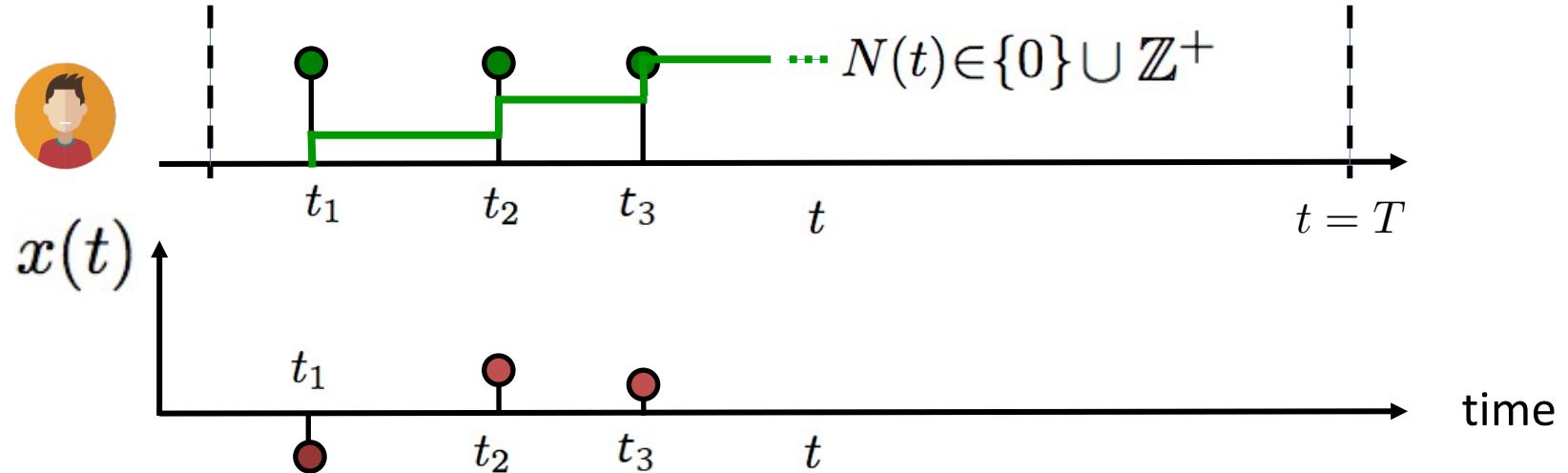
# Marked temporal point processes

Marked temporal point process:

A random process whose realization consists of discrete *marked* events localized in time



# Independent identically distributed marks



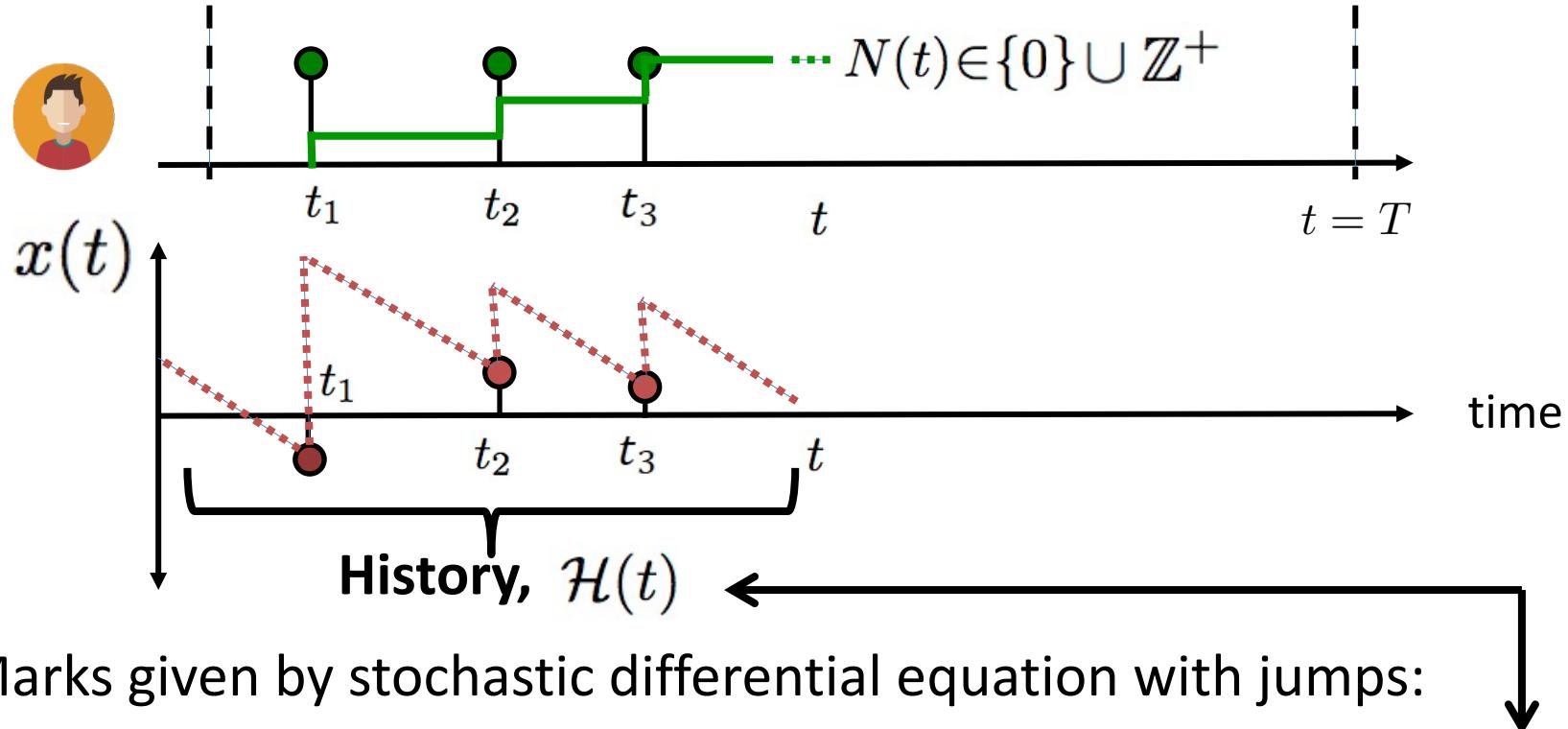
Distribution for the marks:

$$x^*(t_i) \sim p(x)$$

Observations:

1. Marks independent of the temporal dynamics
2. Independent identically distributed (I.I.D.)

# Dependent marks: SDEs with jumps



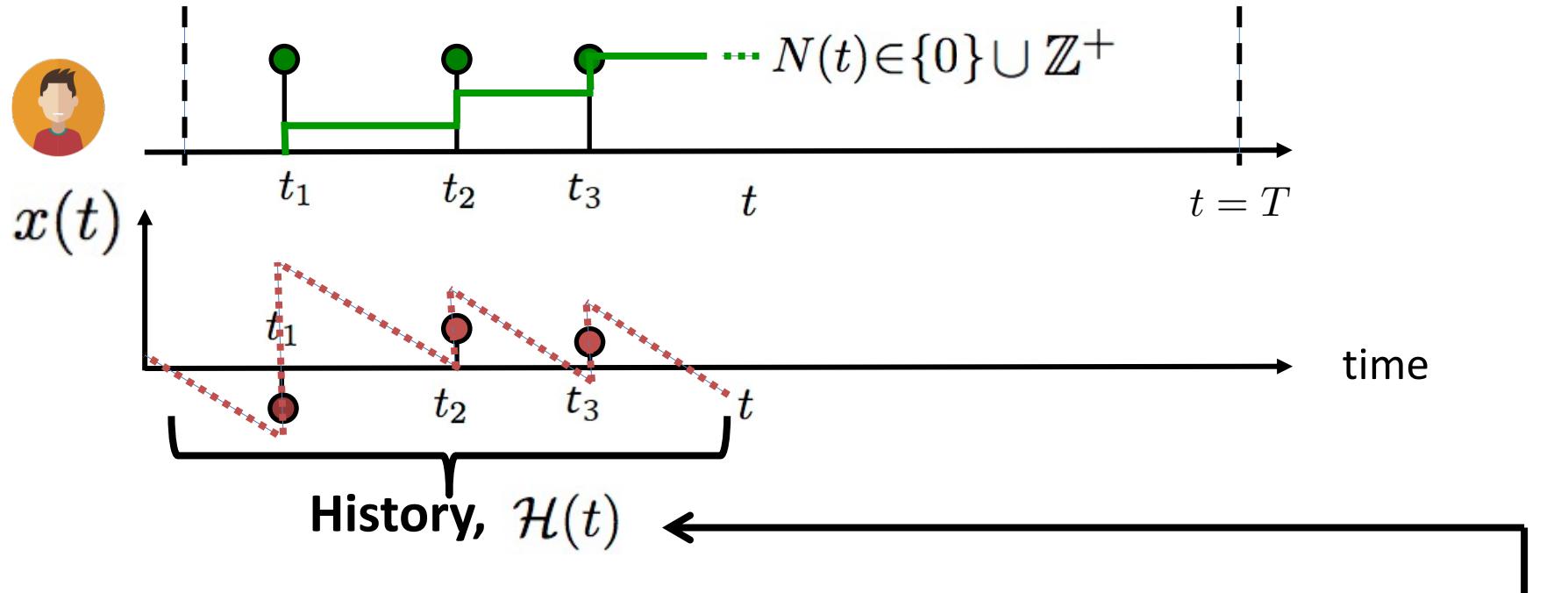
Marks given by stochastic differential equation with jumps:

$$x(t + dt) - x(t) = dx(t) = \underbrace{f(x(t), t)dt}_{\text{Drift}} + \underbrace{h(x(t), t)dN(t)}_{\text{Event influence}}$$

Observations:

1. Marks dependent of the temporal dynamics
2. Defined for all values of  $t$

# Dependent marks: distribution + SDE with jumps



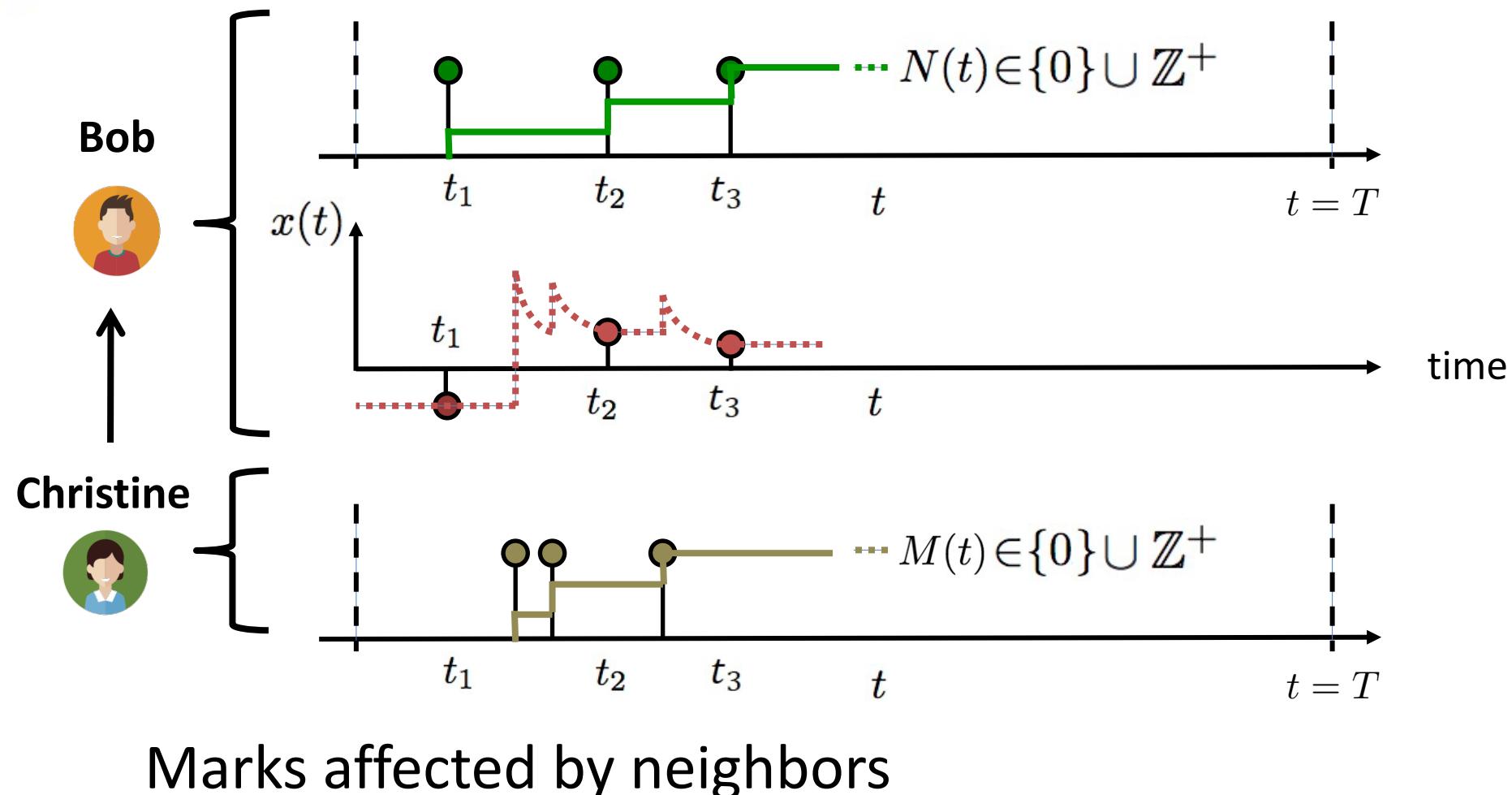
Distribution for the marks:

$$x^*(t_i) \sim p(x^* | x(t)) \Rightarrow dx(t) = \underbrace{f(x(t), t)dt}_{\text{Drift}} + \underbrace{h(x(t), t)dN(t)}_{\text{Event influence}}$$

Observations:

1. Marks dependent on the temporal dynamics
2. Distribution represents additional source of uncertainty

# Mutually exciting + marks



$$dx(t) = \underbrace{f(x(t), t)dt}_{\text{Drift}} + \underbrace{g(x(t), t)dM(t)}_{\text{Neighbor influence}}$$

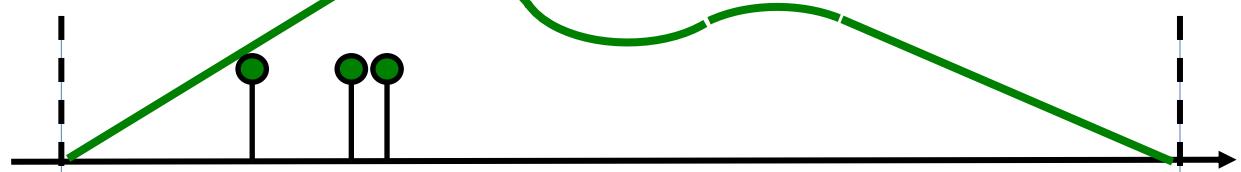


# **Models & Inference: Neural networks for the win**

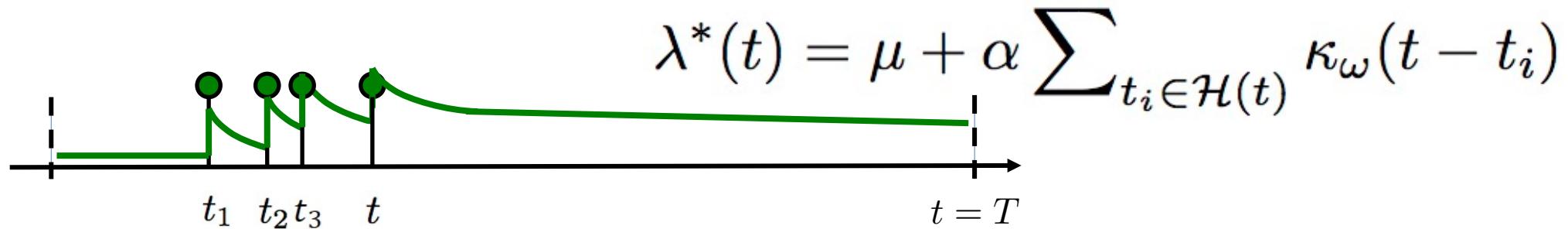
# Hawkes process

We focused on simple temporal dynamics and corresponding intensity functions:

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



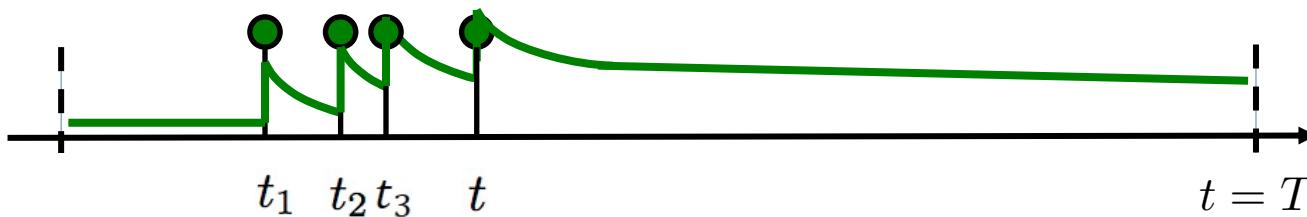
$$\lambda^*(t) = \sum_j \alpha_j k(t - t_j)$$



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

# Triggering Kernels for Hawkes point process

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



Triggering kernel is a linear sum of kernels

$$g_{uu'}(t) = \sum_{d=1}^D a_{uu'}^d g_d(t),$$

Loss function that penalizes non-smooth kernels

$$\mathcal{L}_\alpha(\Theta) = -\mathcal{L}(\Theta) + \alpha \left( \sum_d \mathcal{R}(g_d) + \sum_{u,u',d} (a_{uu'}^d)^2 \right).$$

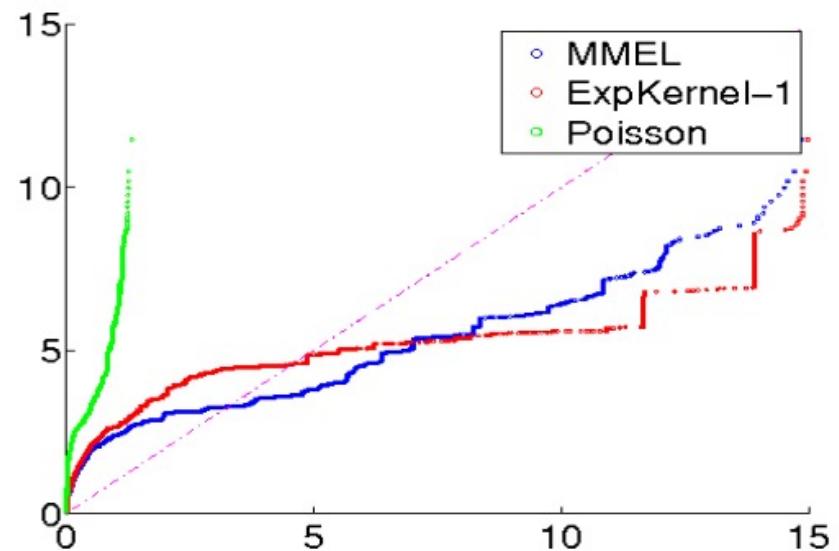
# Optimization of kernels

We define  $Q(\Theta; \Theta^{(k)})$  s.t.

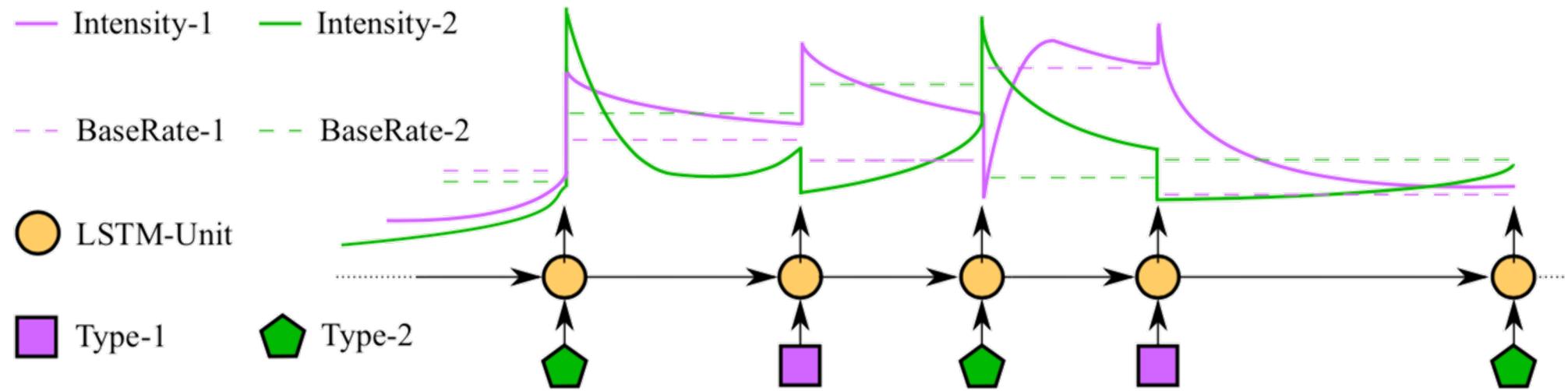
1. For all  $\Theta$  and  $\Theta^{(k)}$ ,  $Q(\Theta; \Theta^{(k)}) \geq \mathcal{L}_\alpha(\Theta)$ .
2.  $Q(\Theta^{(k)}; \Theta^{(k)}) = \mathcal{L}_\alpha(\Theta^{(k)})$ .

The above two properties imply that if  $\Theta^{(k+1)} = \operatorname{argmin}_\Theta Q(\Theta; \Theta^{(k)})$ , we have  $\mathcal{L}_\alpha(\Theta^{(k)}) \geq \mathcal{L}_\alpha(\Theta^{(k+1)})$ .

Separately update kernels and intensities + kernel weights.



# Neural Hawkes process



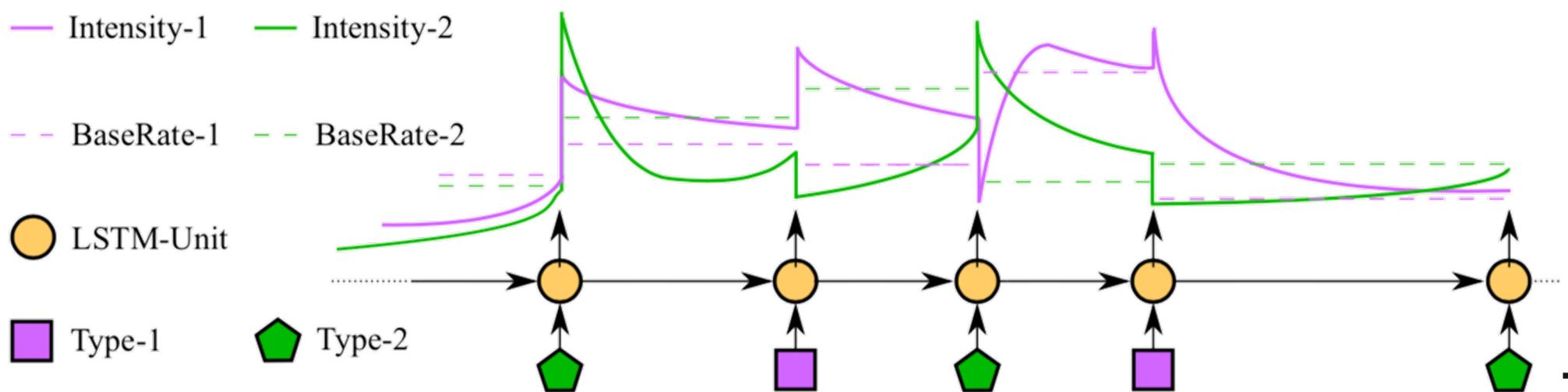
Events affect the intensity of the process in a complex way

- 1) History effect does not need to be additive
- 2) Allows for complex memory effects  
e.g. *delays*

# Neural Hawkes process at event times

$$\mathbf{h}(t) = \text{RNN}(\mathcal{H}(t)) \quad \text{memory via the continuous-time LSTM}$$

$$\lambda_u(t) = f_u(\mathbf{w}_u^\top \mathbf{h}(t)) \quad \text{excitation & inhibition via activation function}$$



# Continuous intensity function for Neural Hawkes

Two cell states allow us to define continuous cell state:

$$\mathbf{c}(t) \stackrel{\text{def}}{=} \bar{\mathbf{c}}_{i+1} + (\mathbf{c}_{i+1} - \bar{\mathbf{c}}_{i+1}) \exp(-\delta_{i+1}(t - t_i)) \text{ for } t \in (t_i, t_{i+1}]$$

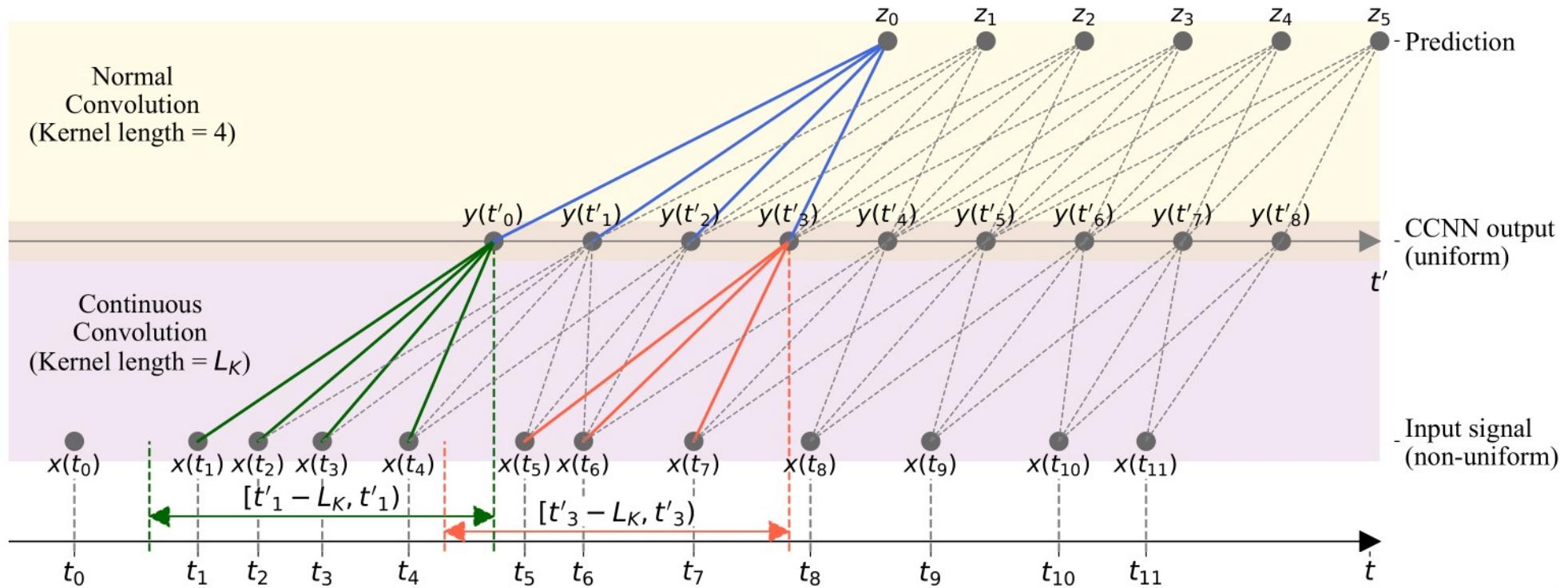
Continuous hidden state is a function of continuous cell state:

$$\mathbf{h}(t) = \mathbf{o}_i \odot (2\sigma(2\mathbf{c}(t)) - 1) \text{ for } t \in (t_{i-1}, t_i]$$

Now we can define the intensity function:

$$\lambda_u(t) = f_u(\mathbf{w}_u^\top \mathbf{h}(t))$$

# CNNs for point processes



After uniform resampling at the first layer we can apply normal convolutions

# Continuous time CNNs

The convolution is defined in the following way for input  $\mathbf{x}$  and kernel  $\psi$ :

$$(\mathbf{x} * \psi)(t) = \sum_{c=1}^{N_{\text{in}}} \int_{\mathbb{R}} x_c(\tau) \psi_c(t - \tau) d\tau.$$

Discretization of the convolution:

$$(\mathbf{x} * \psi)(t) = \sum_{c=1}^{N_{\text{in}}} \sum_{\tau=0}^t x_c(\tau) \psi_c(t - \tau).$$

The kernel is an MLP (Multi-layer perceptron). The activations have *sin* nonlinearities.

$$\mathbf{y} = \text{Sine}(\omega_0[\mathbf{W}\mathbf{x} + \mathbf{b}]).$$

# Intensity function for continuous CNNs

We get intensity by adding state at the last event and weighted time since the last event:

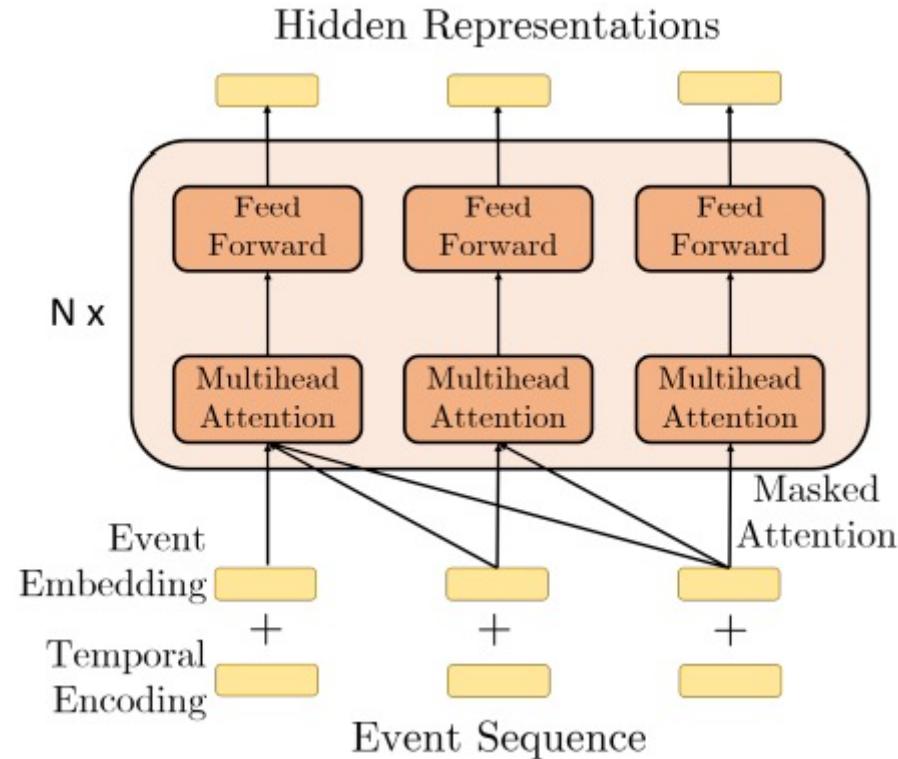
$$\lambda^*(t) = \exp(\mathbf{v}\mathbf{h}_{i-1} + w(t - t_{i-1}) + b),$$

The closed-form expression for conditional intensity is

$$f^*(t) = \exp \left( \mathbf{v}\mathbf{h}_{i-1} + w(t - t_{i-1}) + b + \frac{1}{w} \left( \exp(\mathbf{v}\mathbf{h}_{j-1} + b) - \exp(\mathbf{v}\mathbf{h}_{j-1} + w(t - t_{j-1}) + b) \right) \right).$$

The kernel is an MLP (Multi-layer perceptron).

# Transformers for Neural Hawkes process



We calculate densities and predict next event types:

$$p(t|\mathcal{H}_t) = \lambda(t|\mathcal{H}_t) \exp \left( - \int_{t_j}^t \lambda(\tau|\mathcal{H}_\tau) d\tau \right),$$
$$\hat{t}_{j+1} = \int_{t_j}^\infty t \cdot p(t|\mathcal{H}_t) dt,$$

$$\hat{k}_{j+1} = \operatorname{argmax}_k \frac{\lambda_k(t_{j+1}|\mathcal{H}_{j+1})}{\lambda(t_{j+1}|\mathcal{H}_{j+1})}.$$

Similar idea to calculate intensity for k-th event type in the future:

$$\lambda_k(t|\mathcal{H}_t) = f_k \left( \underbrace{\alpha_k \frac{t - t_j}{t_j}}_{current} + \underbrace{\mathbf{w}_k^\top \mathbf{h}(t_j)}_{history} + \underbrace{b_k}_{base} \right).$$



# **Clustering of event sequences**

# NTPP-MIX framework for the clustering of event sequences

Suitable for any probabilistic model of a point process

EM algorithm with Mean-Field approximation to optimize the cluster parameters and cluster mixing.

After training, an event sequence is input into K NTPPs to get conditional probabilities.

Adding corresponding expectation of mixing coefficients, we get the approximated posterior of cluster assignment.

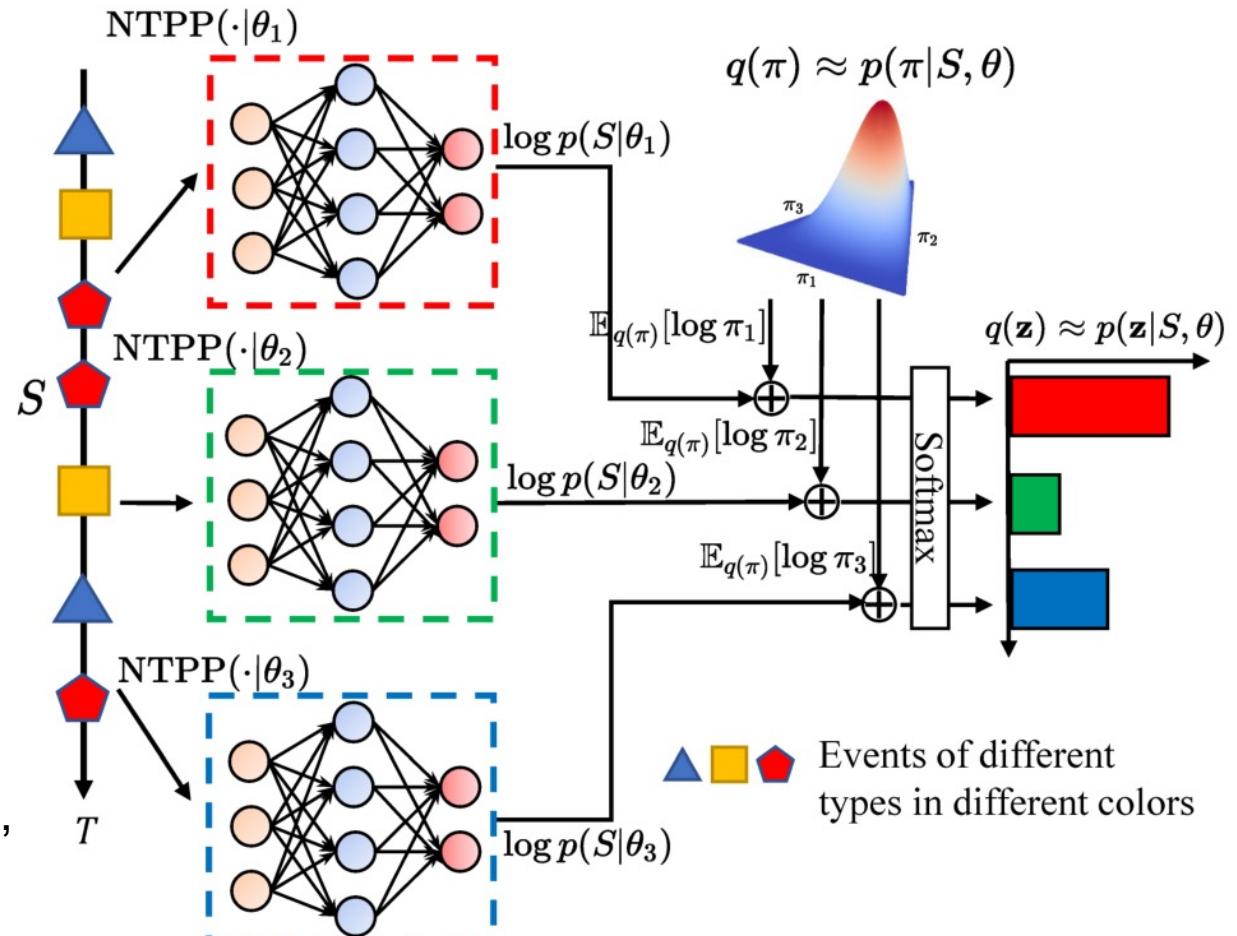
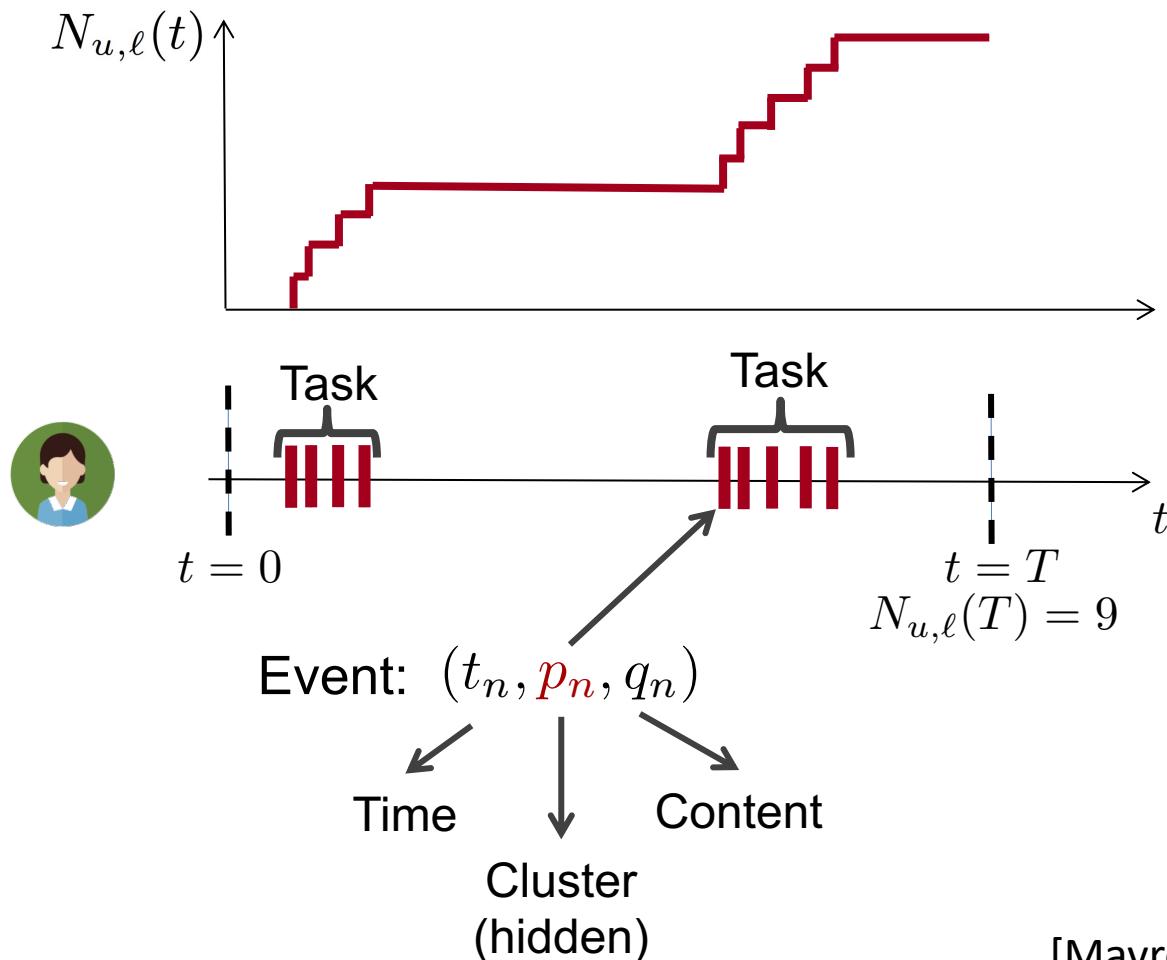


Figure: Inference, example for  $K = 3$

# Events representation

We represent the events using **marked temporal point processes**:



# Conclusions

# Conclusions

- Temporal point processes is a natural model for event sequences data
- Transformers help for incorporation of the whole history