

# Learning from event sequences



**Alexey Zaytsev**

Assistant professor, Skoltech

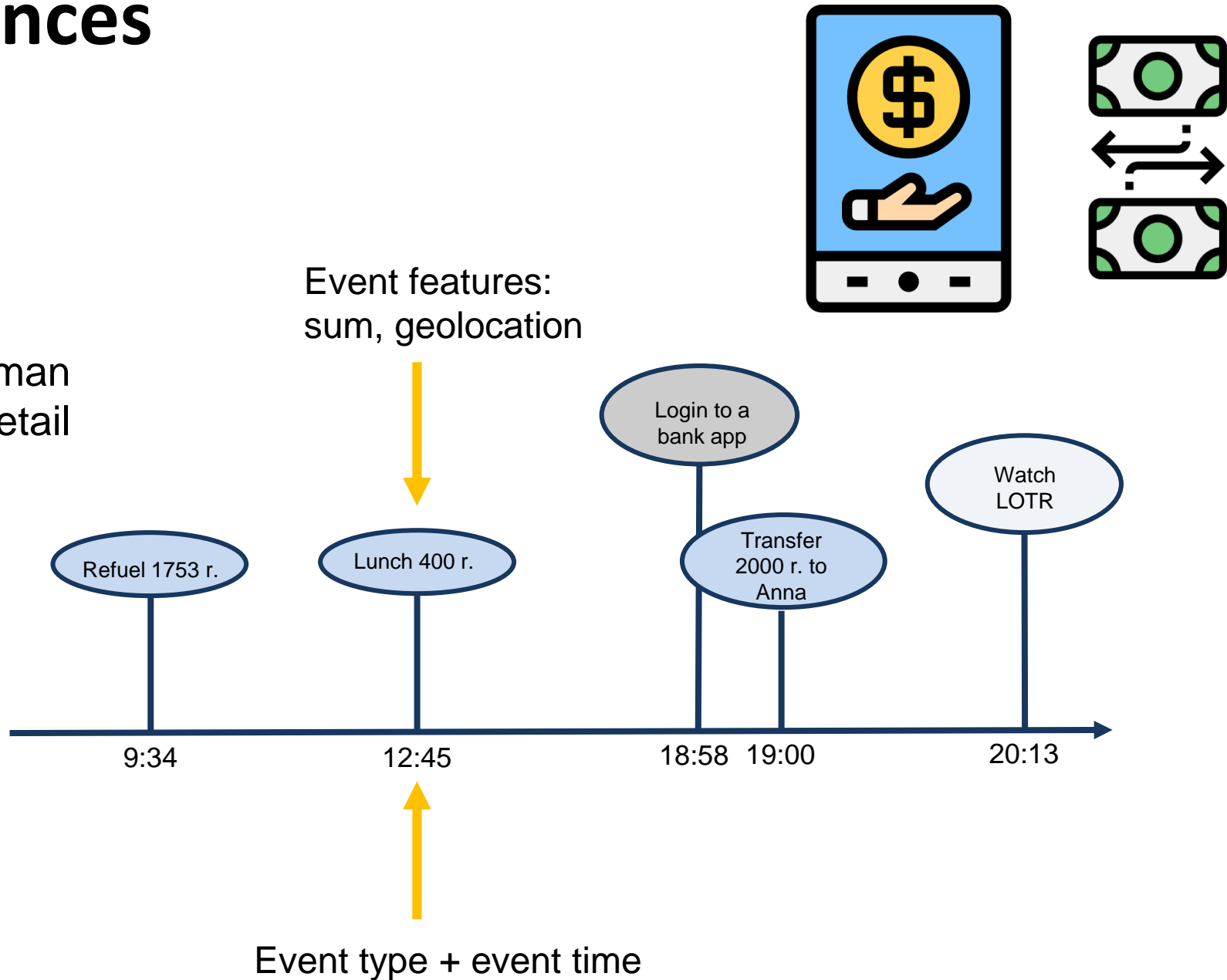


# **Temporal Point Processes (TPPs): applied problems**

# Example: financial transactions are event sequences



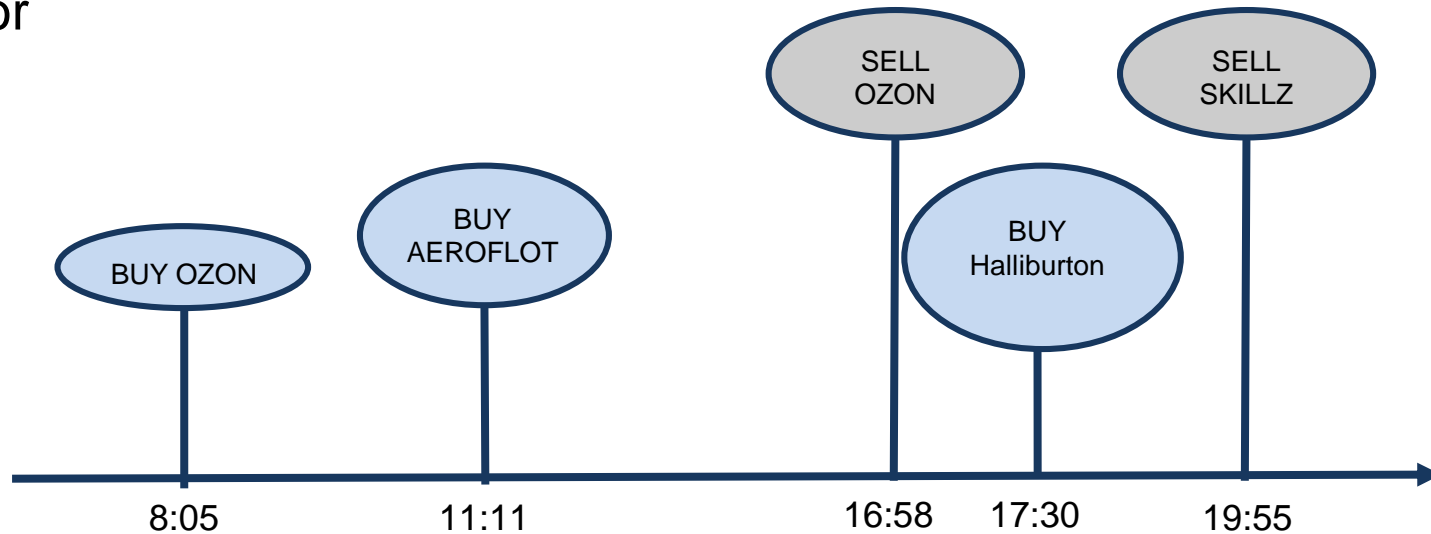
Alex, 33, man  
works in retail



# Example: operations in markets are event sequences



An investor



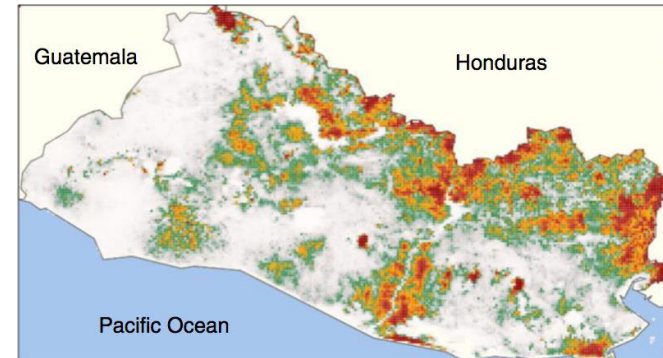
# Many discrete *events* in continuous time



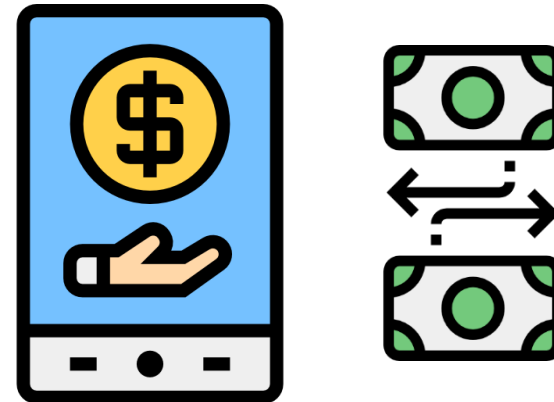
Online actions



Financial trading



Disease dynamics



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

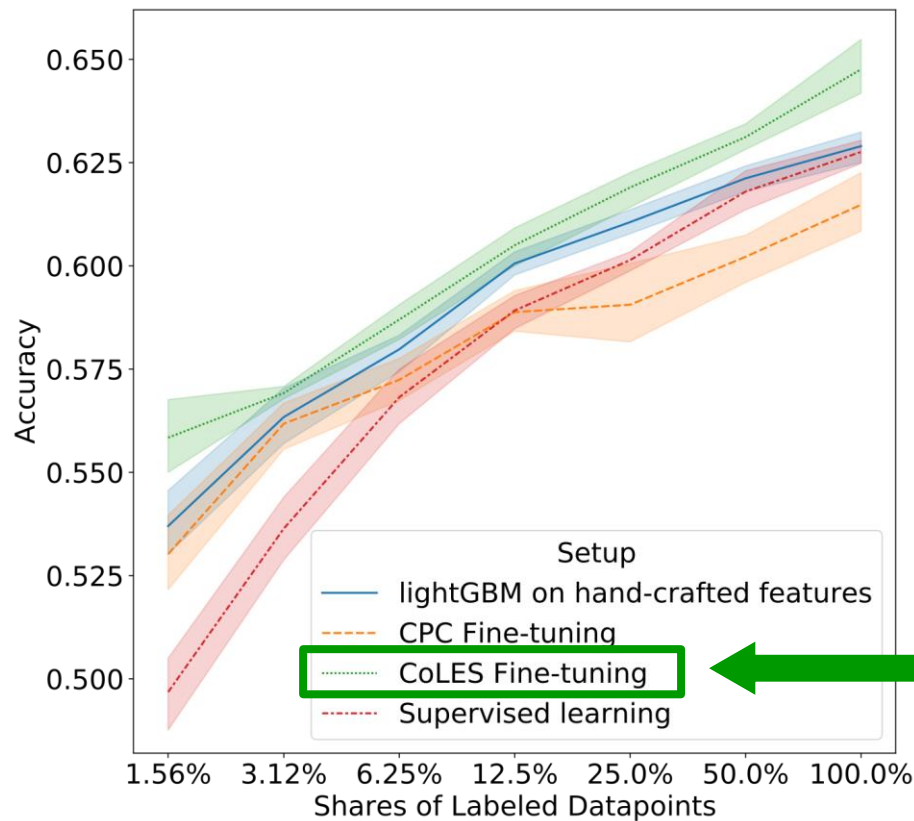
SLOW

...in a wide range of temporal scales.

# Embeddings idea

Financial transactions uncover fundamental information about a customer.

Learn universal embeddings from event sequences that can be used to solve various task in bank.



Embeddings by Sber

# Applications of event sequences in Sber

Embeddings are used in more than 50 different models in the bank

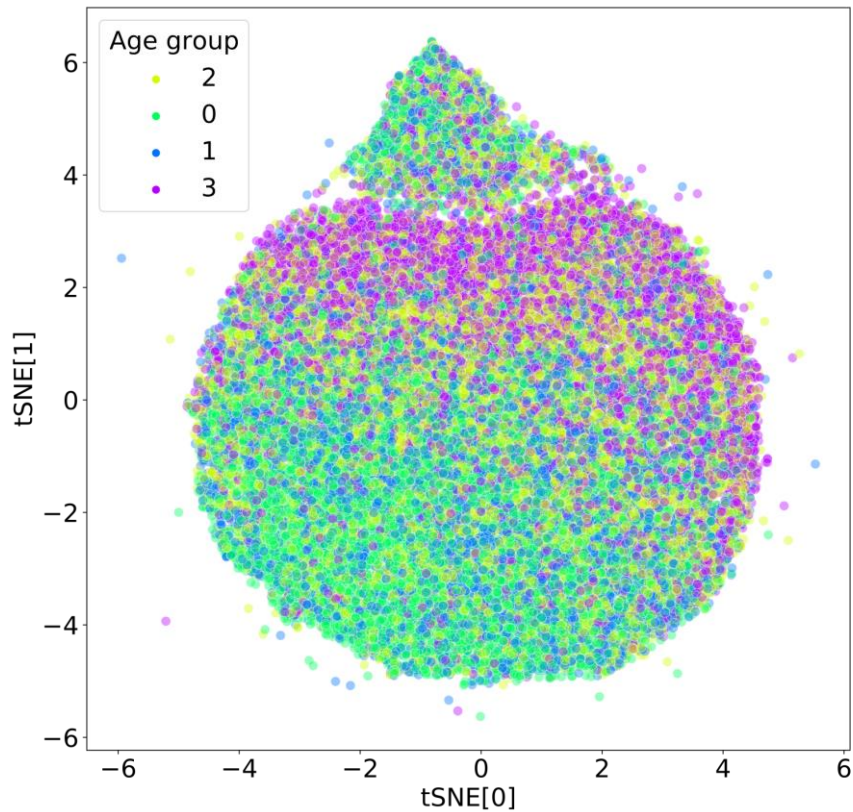
The reason: the quality is superior

Company money transfer history embeddings		Card transaction history embeddings	
Credit product look-alike	+6 Gini	Depositors churn	+5 Gini
Corporate insurance look-alike	+28 Gini	Remote medicine look-alike	+4 Gini
Holding structure prediction	+6 Gini	Movie recommendations cold start	+12,7% NDCG
Scoring for small businesses	+8 Gini	Retail credit scoring	+4 Gini

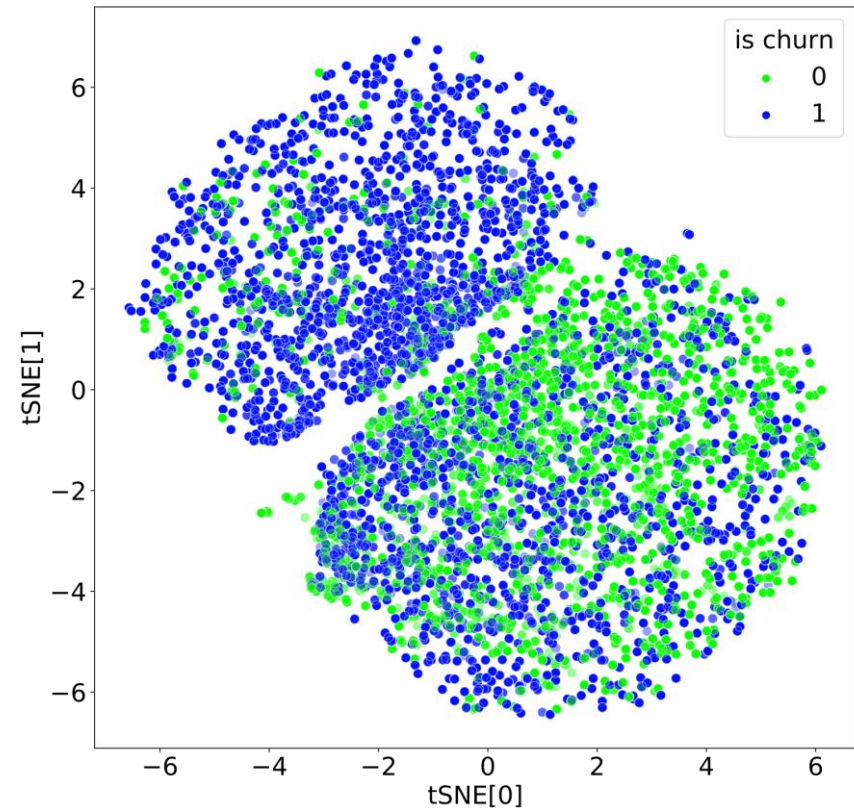


# Let's look at embeddings

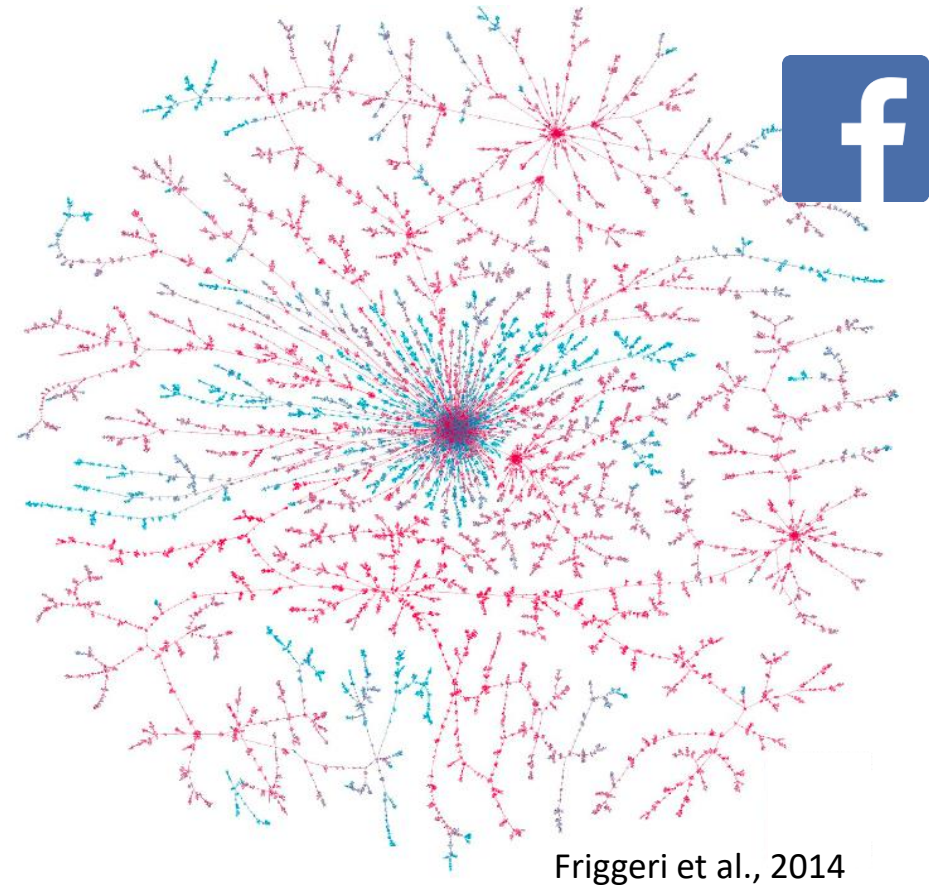
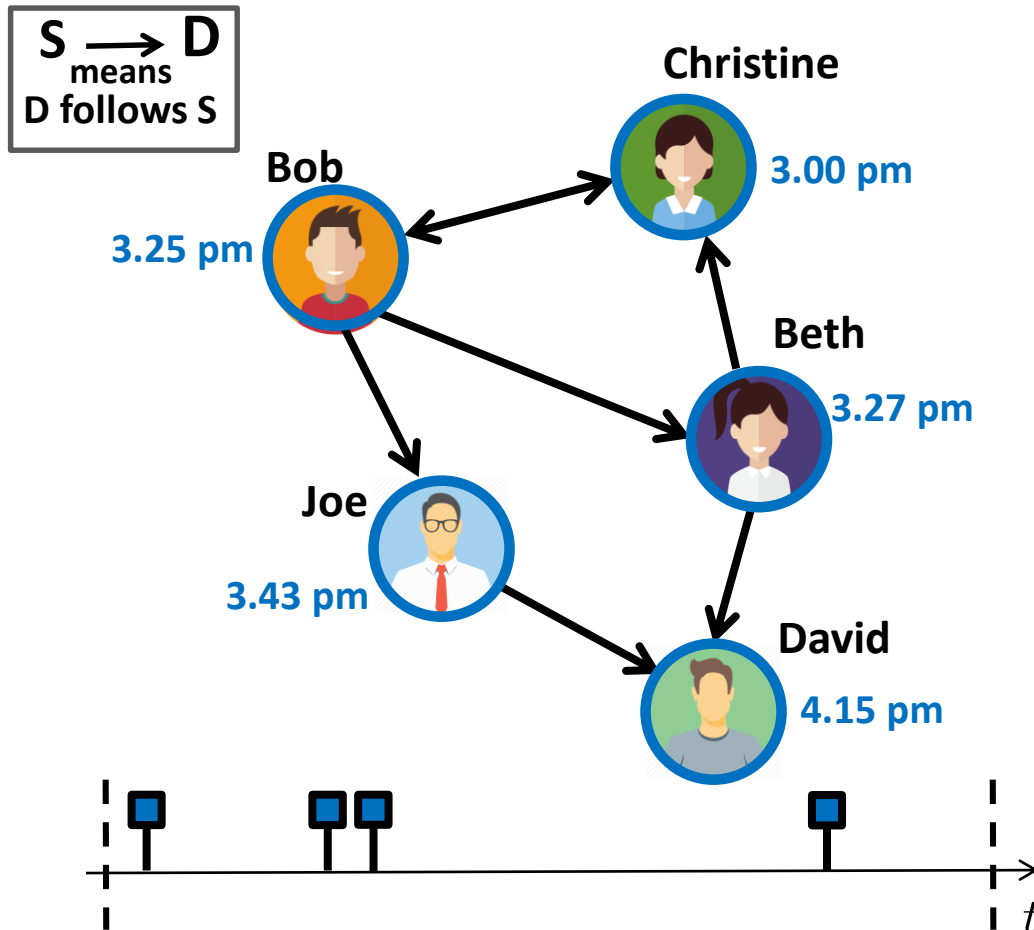
Color shows age group



Color shows churn fact



# Example I: Information propagation



They can have an impact  
in the off-line world

**theguardian**

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

# Example 2: 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 August 17, 1961), is the current President of the United States. He was president of the Harvard Law Review, a civil rights attorney and taught at the University of Chicago Law School, representing the 13th District of Illinois in the United States House of Representatives.

**Barack Obama: Revision history**

03:41, 28 November 2016 Ranze (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) ... (E)  
07:03, 27 November 2016 Saiph121 (talk | contribs) ... (301,192 bytes) (+25) ... (E)

03:21, 20 September 2016

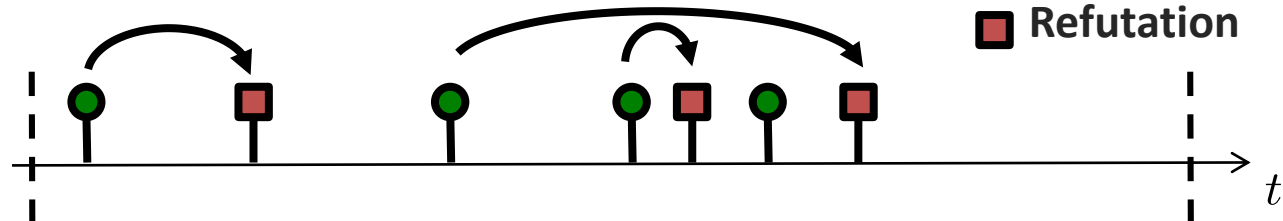
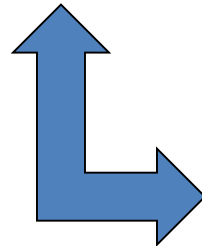
is a **Kenyan** politician



possible vandalism by **MLM2016**

is an American politician

● Addition  
■ Refutation



Moving to Australia Working in Australia Study abroad in Australia +4

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

Answer Request Follow 109 Comment Share 9 Downvote

I have studied, worked and lived in Australia as an International student, business owner and a citizen.

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

Hope it helps! :)

Possible Challenges

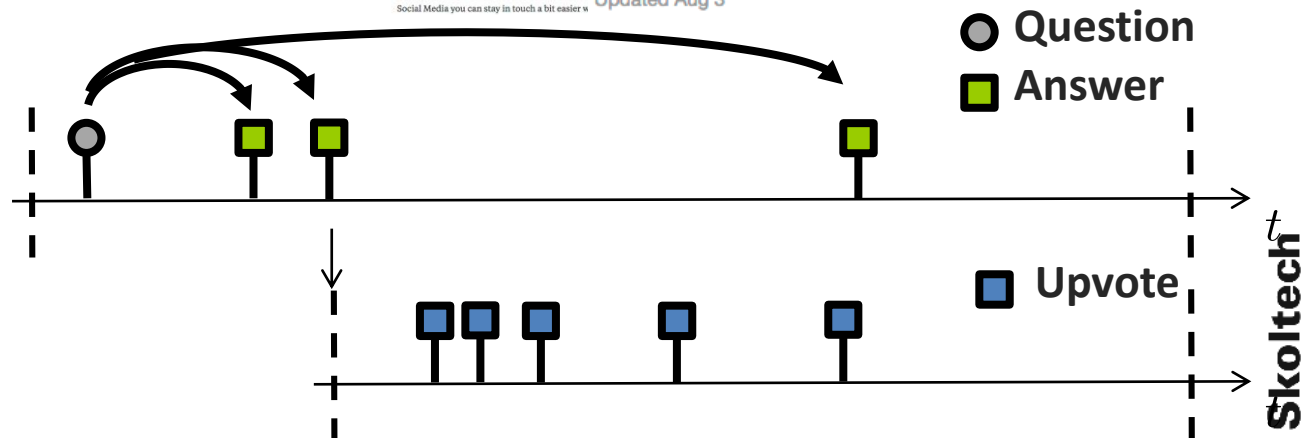
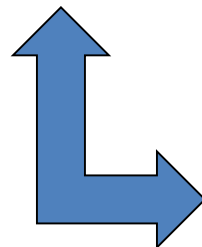
- Language problem for those who don't speak English
- Not having your family and friends around could be a bit lonely
- Society is more and more connected and thanks to Social Media you can stay in touch a bit easier

Upvote 150



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

Updated Aug 3

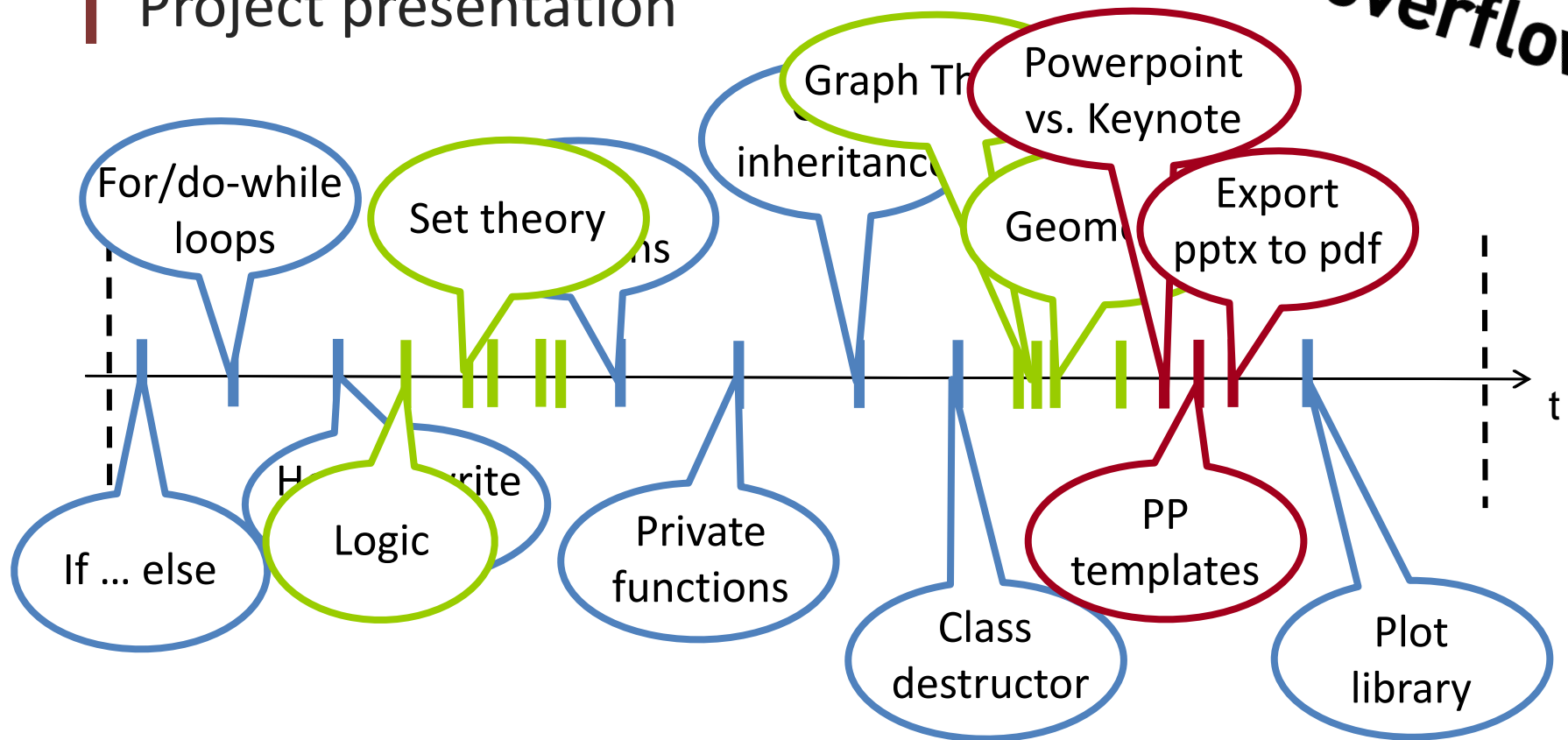


# Example 3: development

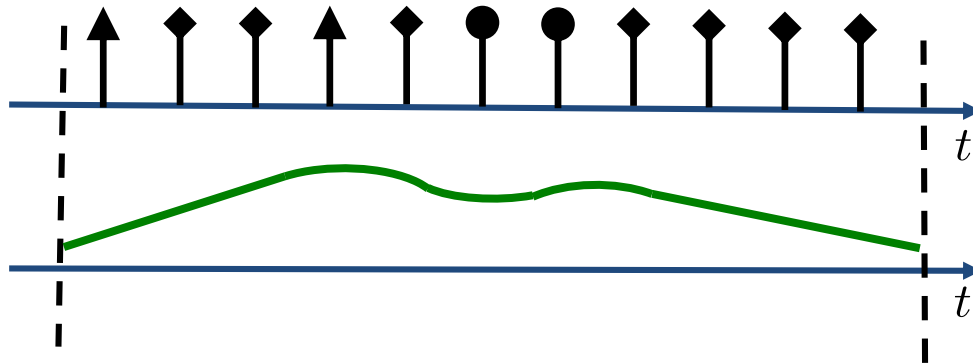


## 1st year computer science student

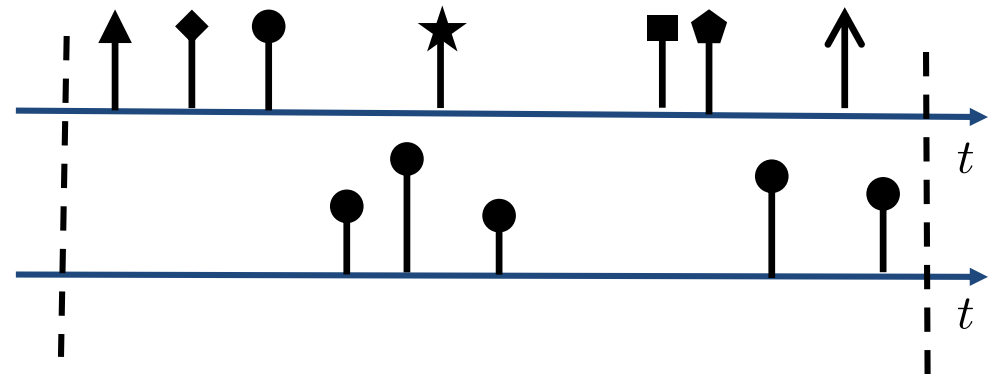
- Introduction to programming
- Discrete math
- Project presentation



# Aren't these event traces just time series?

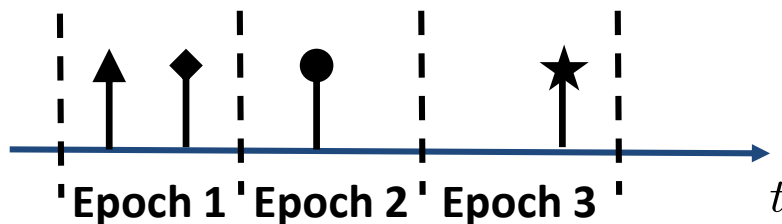


**Discrete and continuous times series**



**Discrete events in continuous time**

**What about aggregating events in *epochs*?**



How long is each epoch?

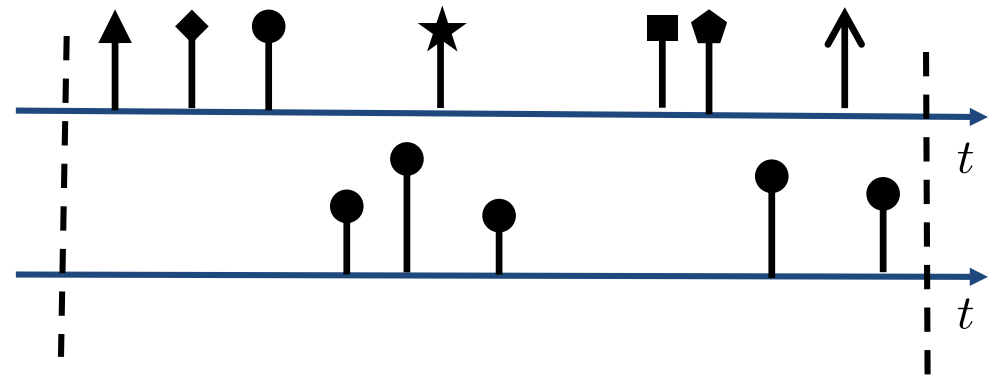
How to aggregate events per epoch?

What if no event in one epoch?

What about time-related queries?

# Problems for event sequences

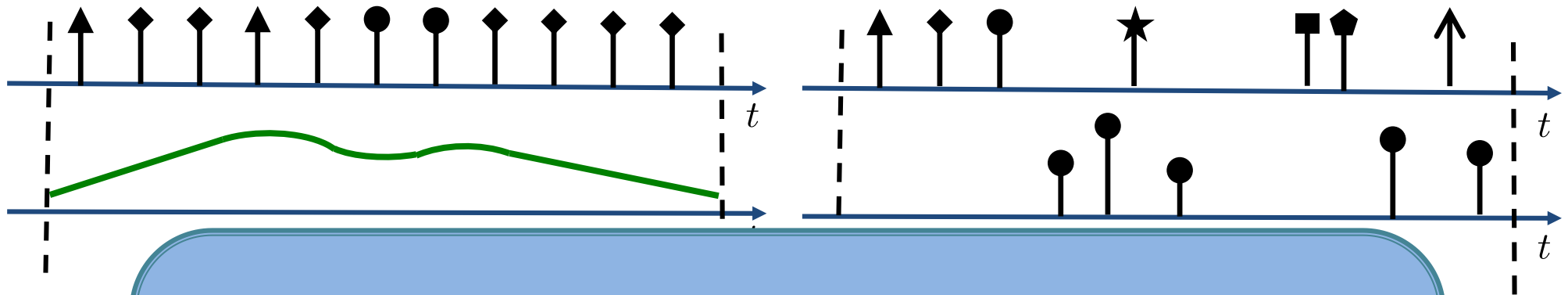
- Compact description of data: models
- Interpretation
- Forecasting/Prediction
- Control
- Hypothesis testing
- Simulation



**Discrete events in  
continuous time**



# Aren't these event traces just time series?



Dis

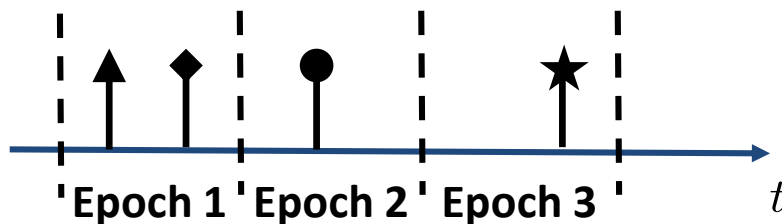
W

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

epoch?

What if no event in one epoch?

What about time-related queries?

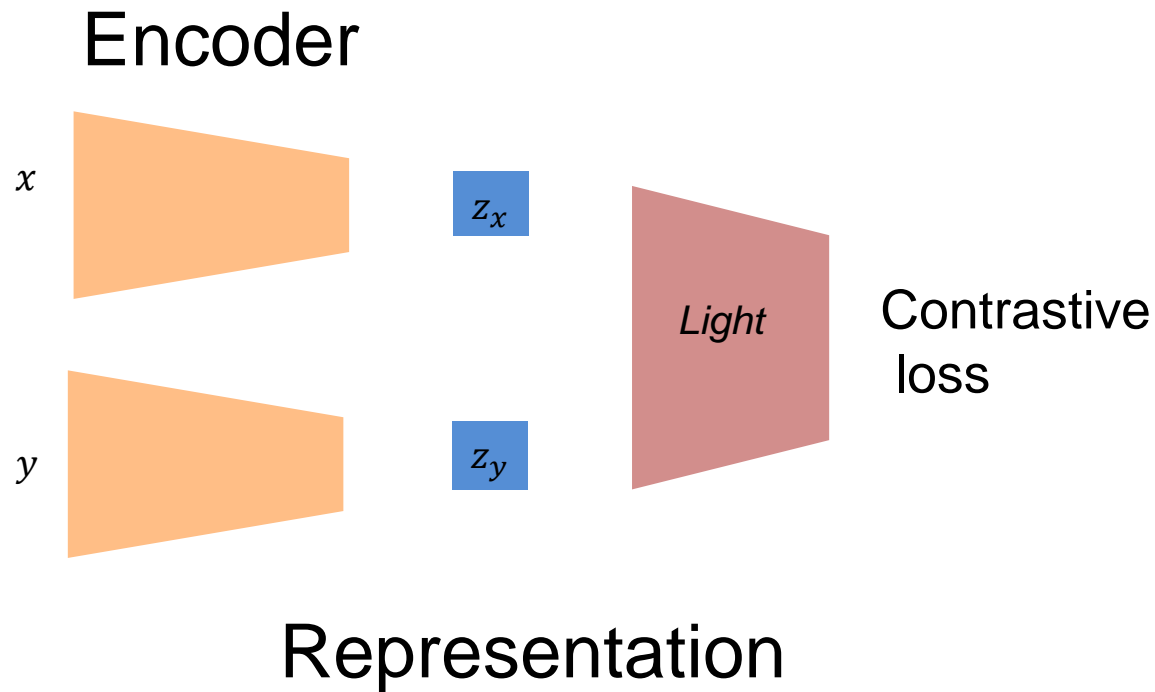


# **Learning universal embeddings**

**Contrastive approach**



# Contrastive self-supervised learning

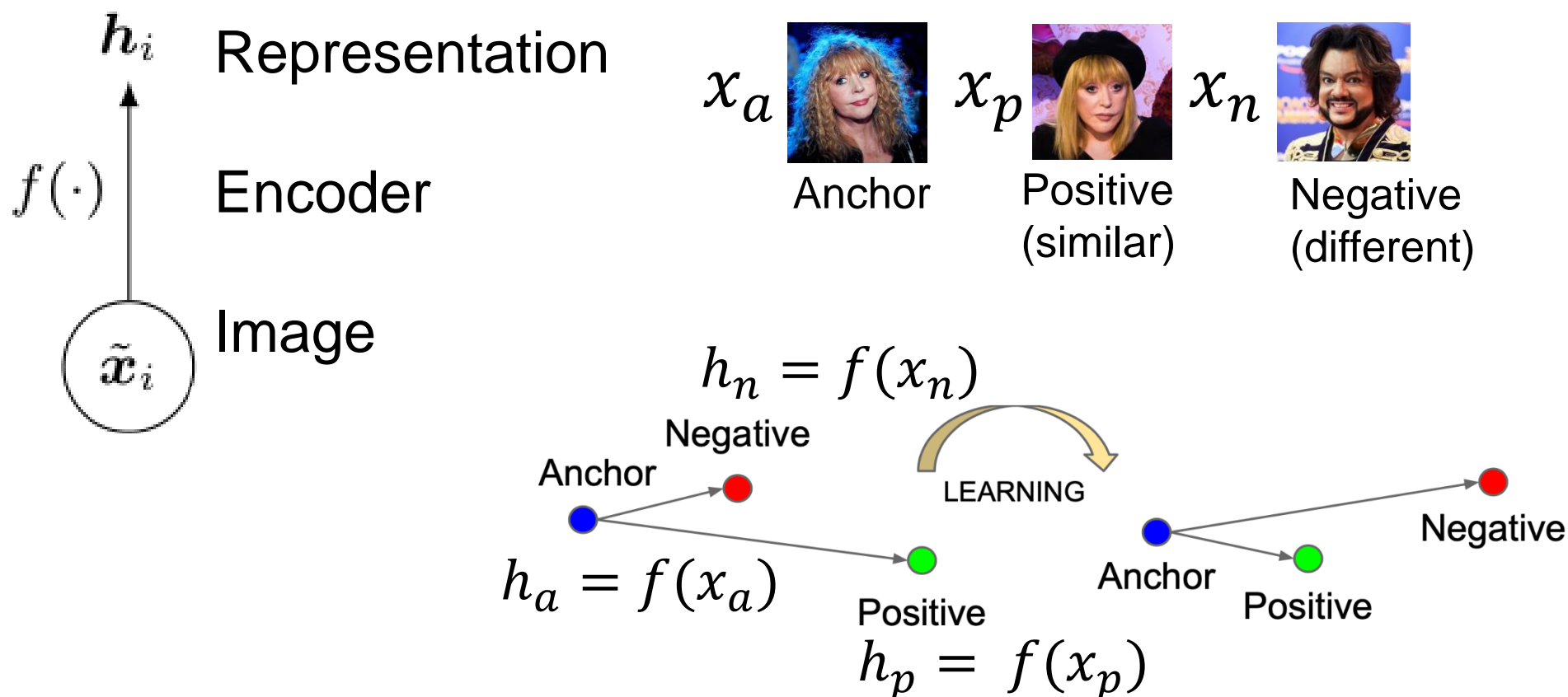


17

Two models: Encoder & Discriminator.  
Encoder produces representations.

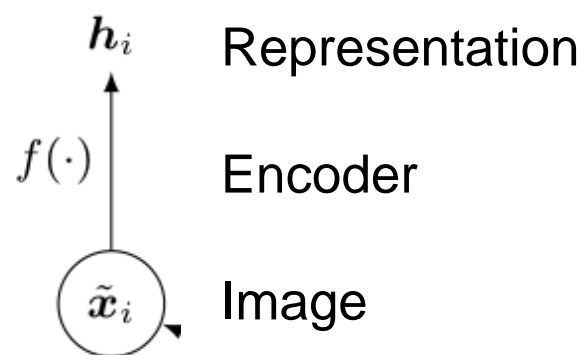
# What if we have labels...

## Contrastive learning idea for supervised learning



# What if we don't have labels...

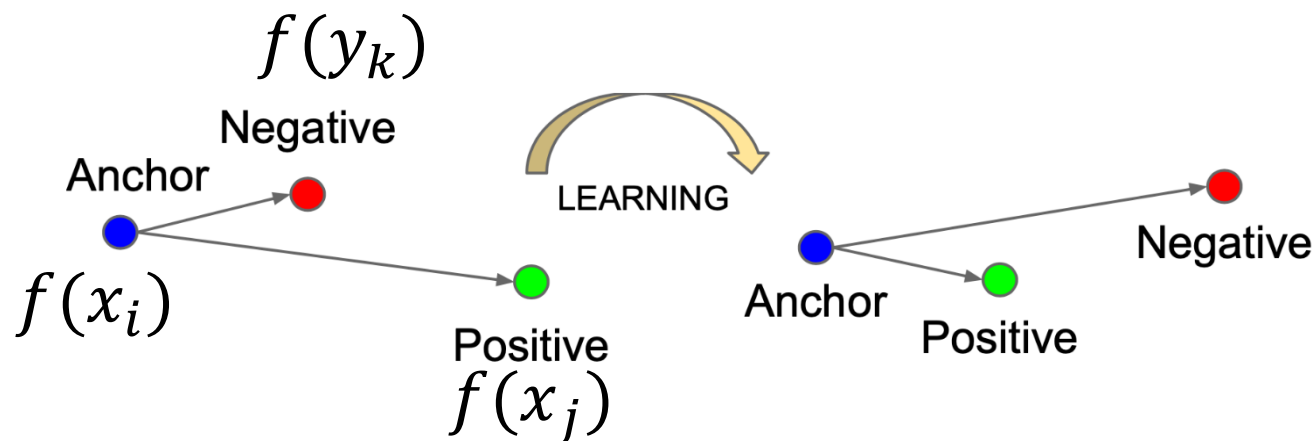
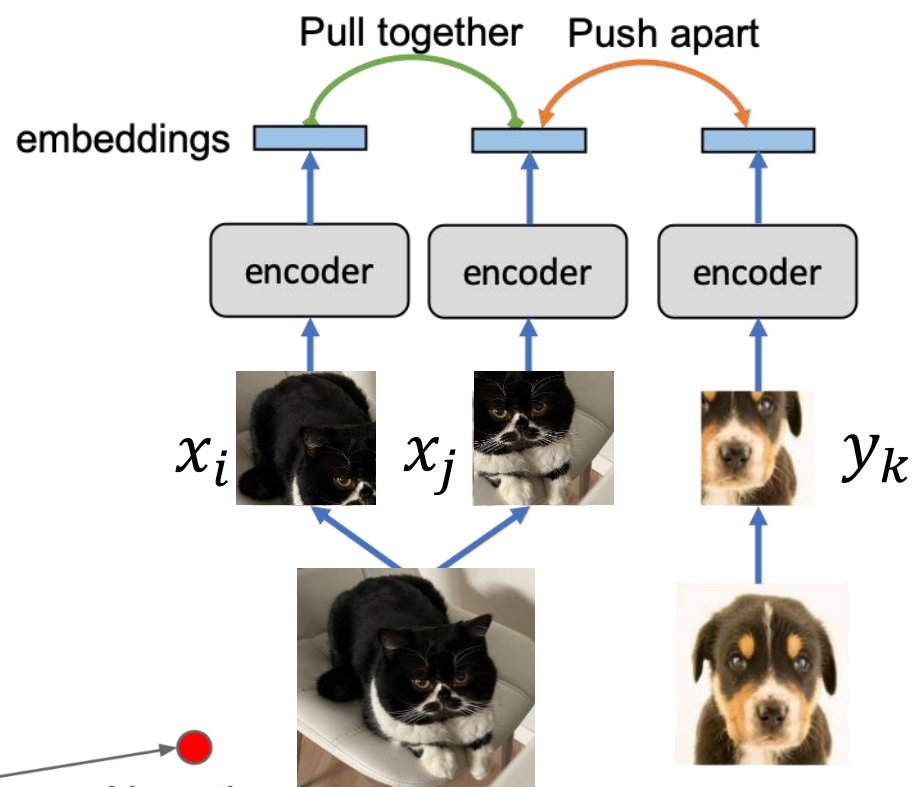
## Generate views!



$x_i$  – one view

$x_j$  – another view

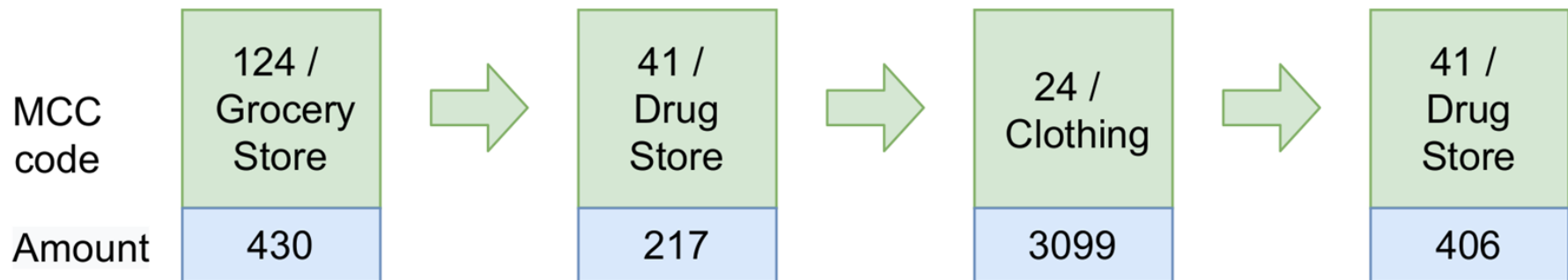
$y_k$  – negative example



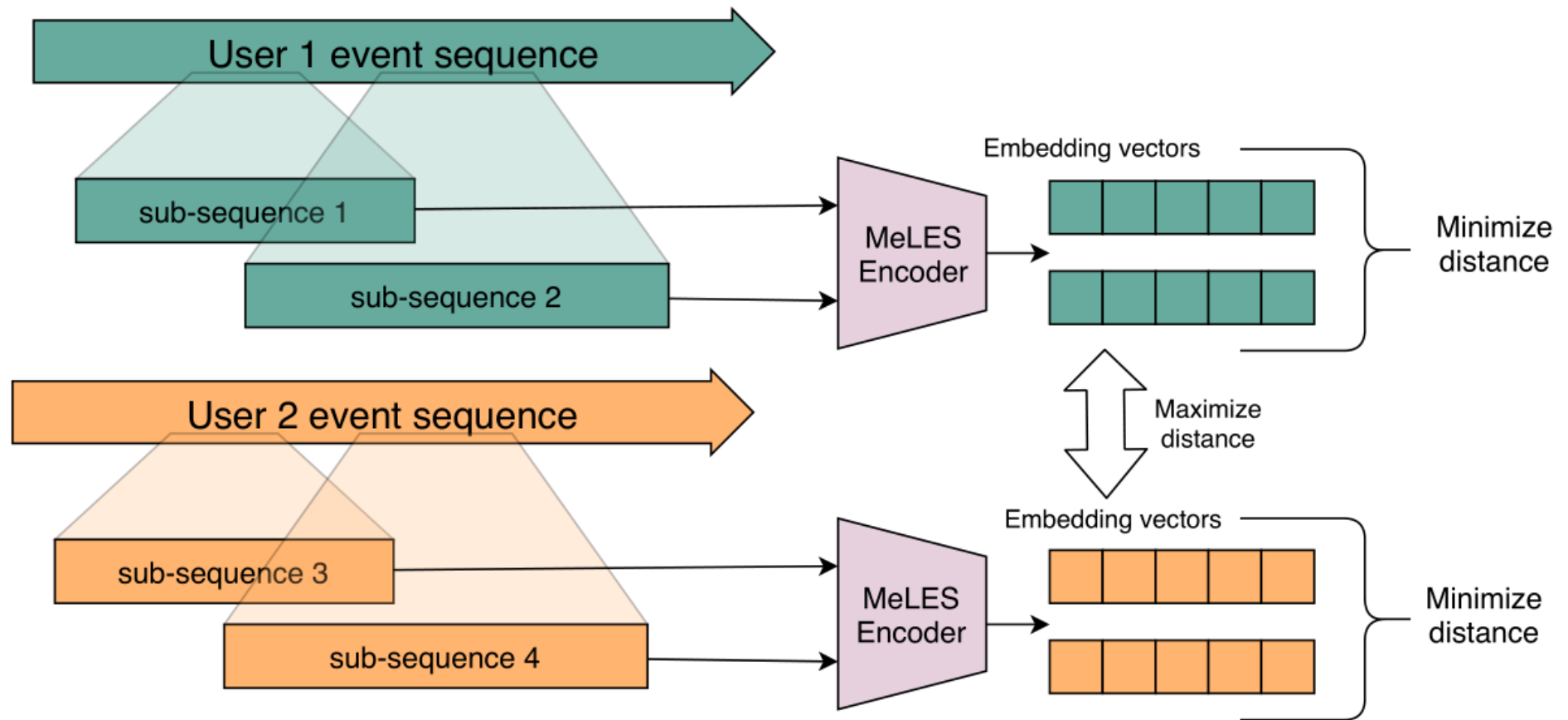
# Discrete sequential transactional data

Transaction records data sequence includes:

- MCC (Merchant Category Codes)
- Purchase amount
- Time values
- Transaction location
- ...



# MeLES contrastive learning



<https://github.com/dl11lb/pytorch-lifestream>

Babaev, Dmitrii, et al. "Event sequence metric learning" KDD. 2022.

# Usage of embeddings

Obtain embeddings from our event processing model.

Train a model on top of embeddings:

- SVM
- Gradient boosting
- Logistic regression

Doesn't need to access to an initial model

Alternatives:

- Fine-tuning all model: longer, required for weaker models)
- Use kNN on top of embeddings: requires stronger model

# Transformers or not?

Table 4: Comparison of encoder types

Encoder type	Age, Accuracy $\pm 95\%$	Gender, AUROC $\pm 95\%$
LSTM	$0.620 \pm 0.003$	$0.870 \pm 0.005$
GRU	<b><math>0.639 \pm 0.006</math></b>	<b><math>0.871 \pm 0.004</math></b>
Transformer	$0.621 \pm 0.001$	$0.848 \pm 0.002$

Table 5: Comparison of metric learning losses

Loss type	Age, Accuracy $\pm 95\%$	Gender, AUROC $\pm 95\%$
Contrastive loss	$0.639 \pm 0.006$	<b><math>0.871 \pm 0.003</math></b>
Binomial deviance loss	$0.535 \pm 0.005$	$0.853 \pm 0.005$
Histogram loss	<b><math>0.642 \pm 0.002</math></b>	$0.851 \pm 0.004$
Margin loss	$0.631 \pm 0.003$	<b><math>0.871 \pm 0.004</math></b>
Triplet loss	$0.610 \pm 0.006$	$0.855 \pm 0.003$

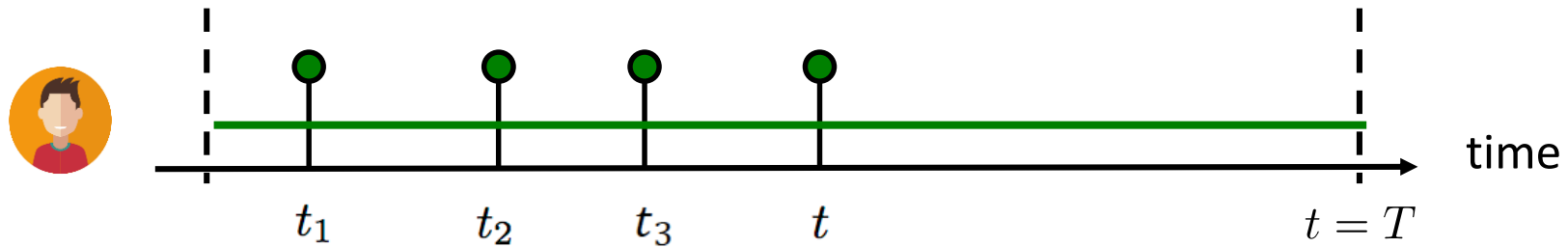
# **Temporal Point Processes (TPPs): Introduction**

**Poisson process and intensity  
function**



# Poisson process

Coarse approximation of many real-life processes



Intensity is the expected number of events per unit time.

Intensity of a Poisson process is constant:

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

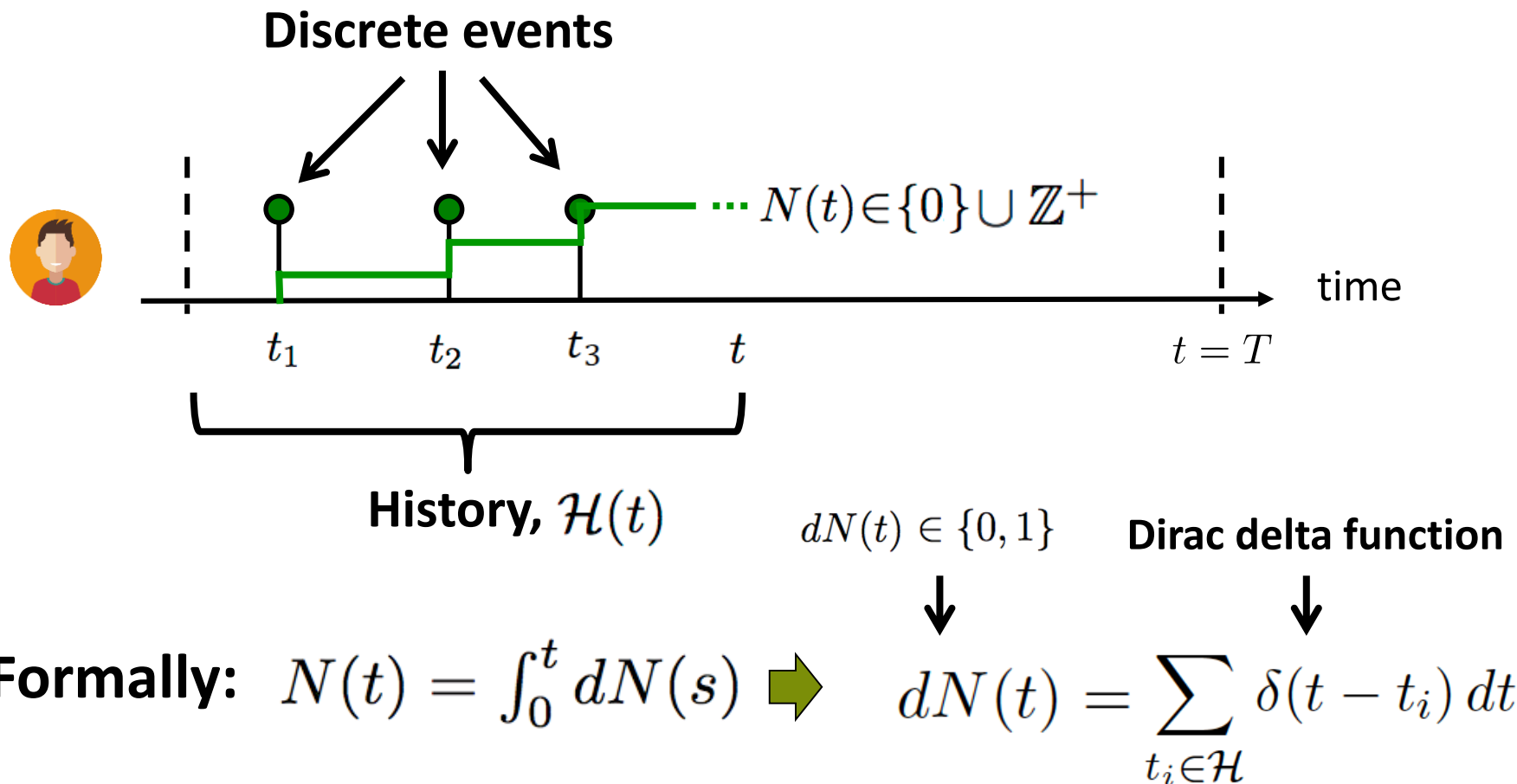
Observations:

1. Intensity independent of history
2. Uniformly random occurrence

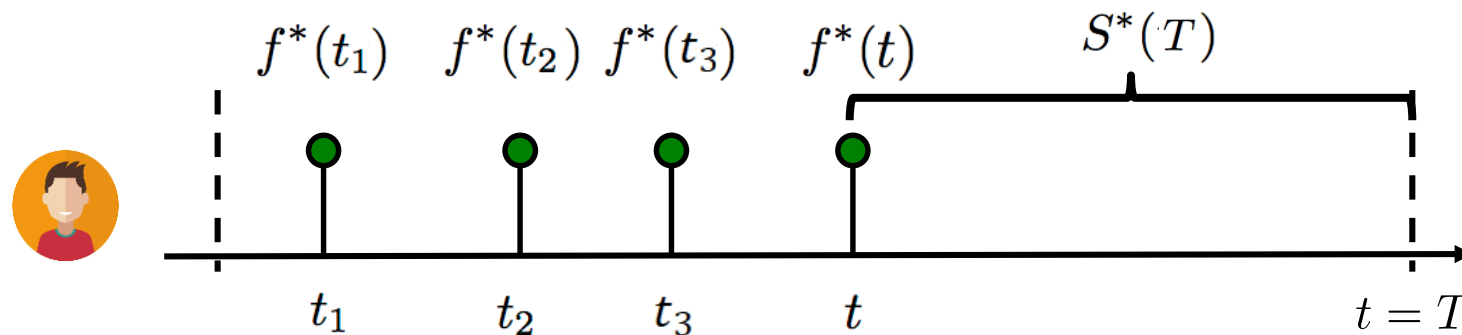
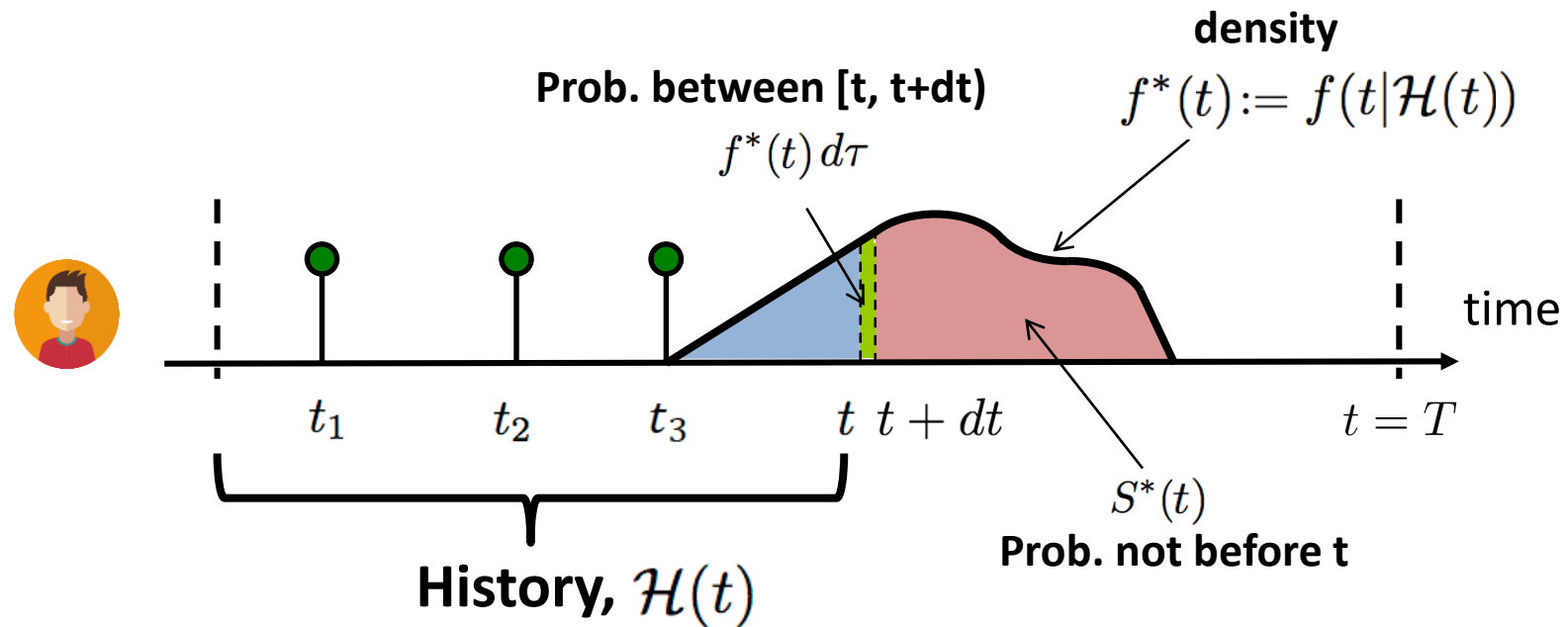
# Temporal point processes

## Temporal point process:

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

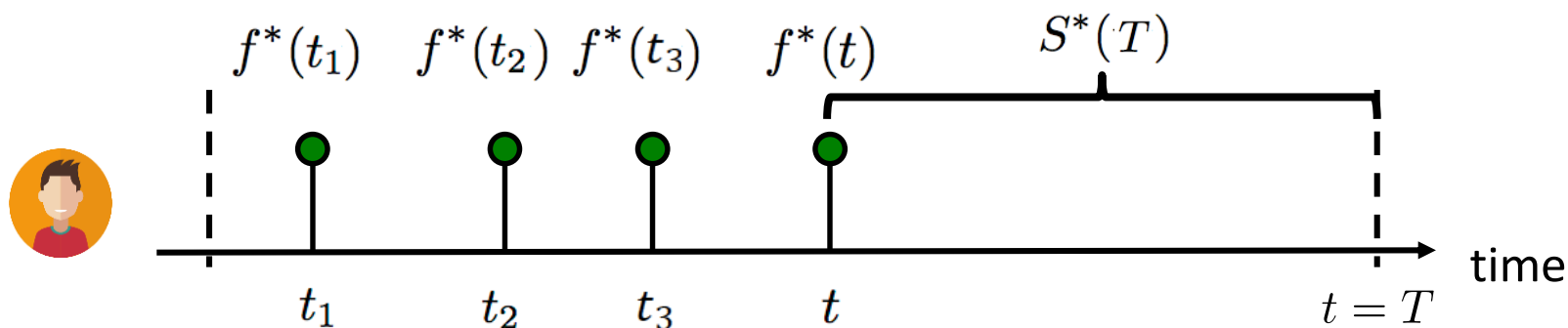


# Model time as a random variable



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

# Problems of density parametrization (I)

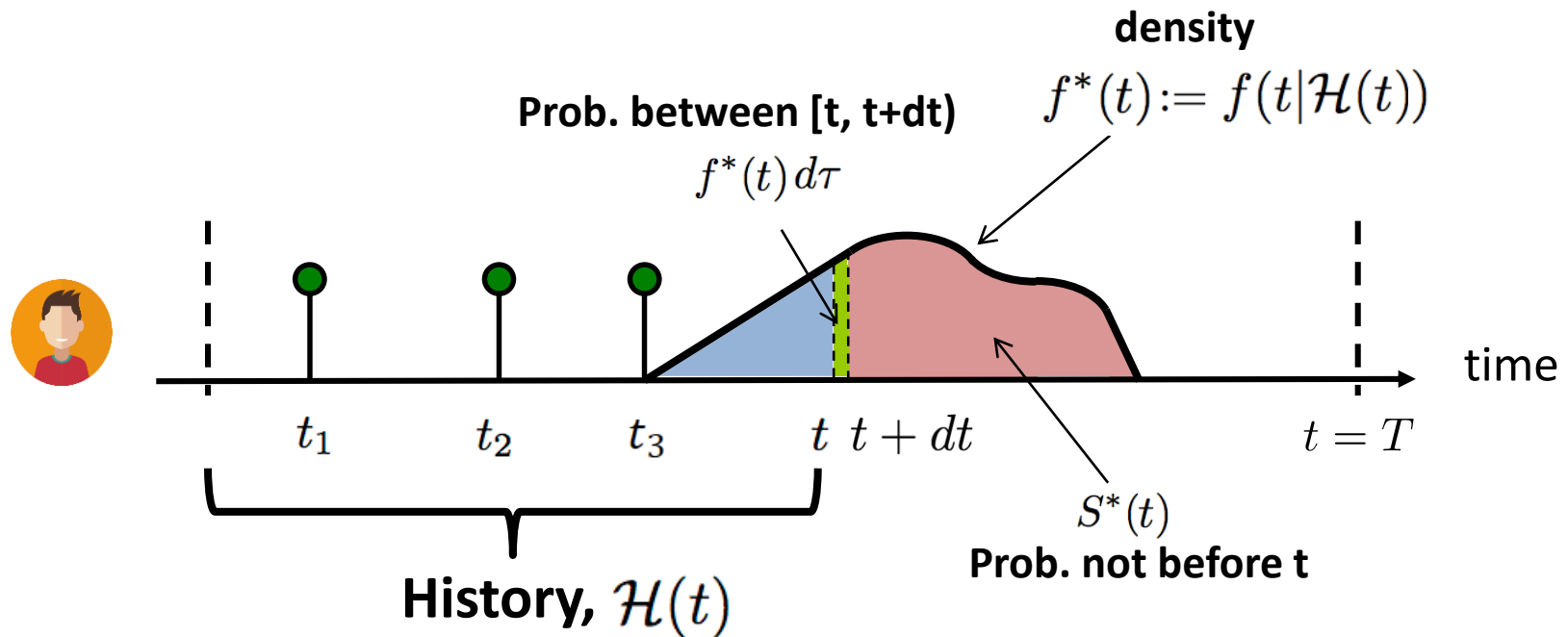


$$\begin{array}{ccccccc}
 f^*(t_1) & f^*(t_2) & f^*(t_3) & f^*(t) & S^*(T) & & \\
 \nearrow & \nearrow & \uparrow & \nwarrow & \nwarrow & & \\
 \frac{\exp\langle w, \psi^*(t_1) \rangle}{Z} & \frac{\exp\langle w, \psi^*(t_2) \rangle}{Z} & \frac{\exp\langle w, \psi^*(t_3) \rangle}{Z} & \frac{\exp\langle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

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



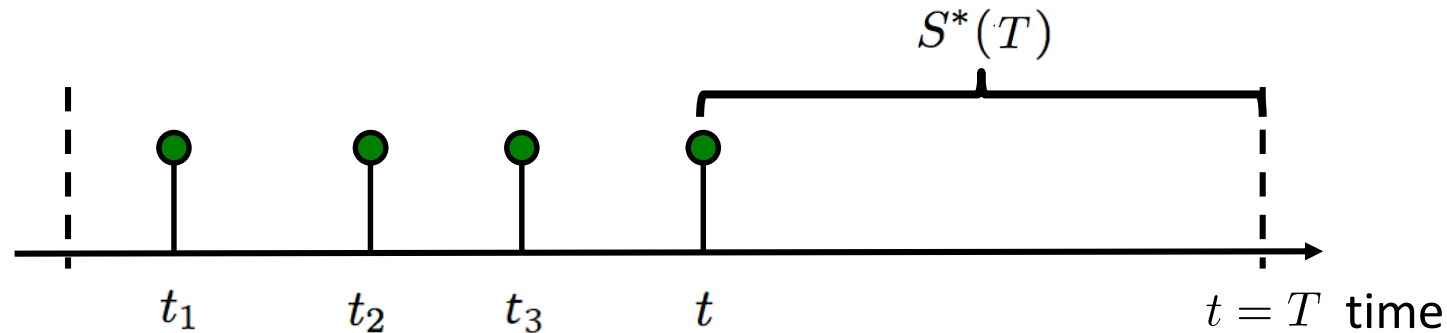
**Intensity:**

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

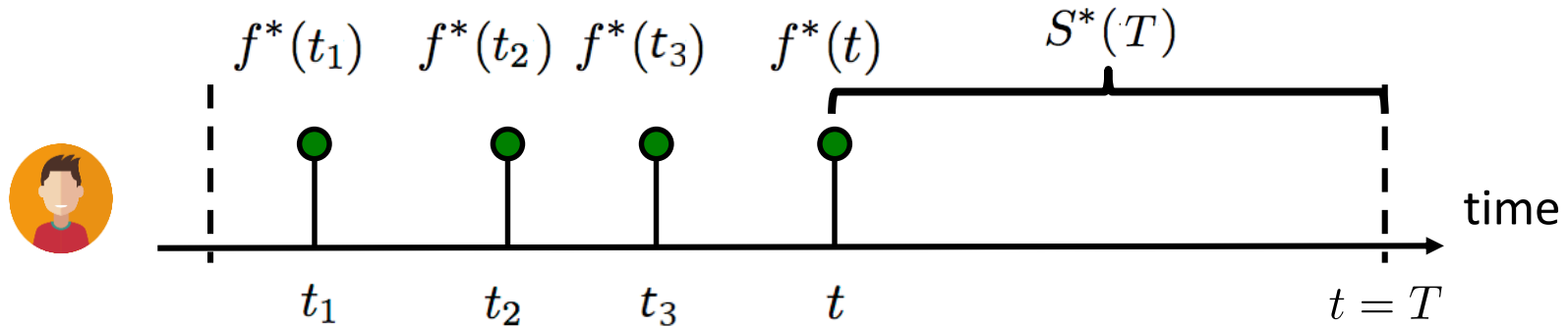
# Log likelihood via intensity and density



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

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

# Advantages of intensity parametrization (I)



$$\lambda^*(t_1) \lambda^*(t_2) \lambda^*(t_3) \lambda^*(t) \exp \left( - \int_0^T \lambda^*(\tau) d\tau \right)$$

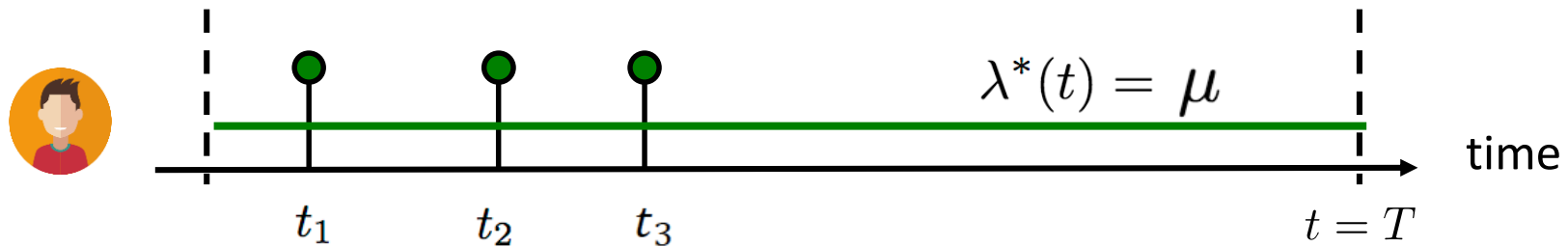
Arrows point from the following expressions to the corresponding terms in the equation above:

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

**Suitable for model design and interpretable:**

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

# Fitting & sampling from a Poisson



Fitting by maximum likelihood:

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

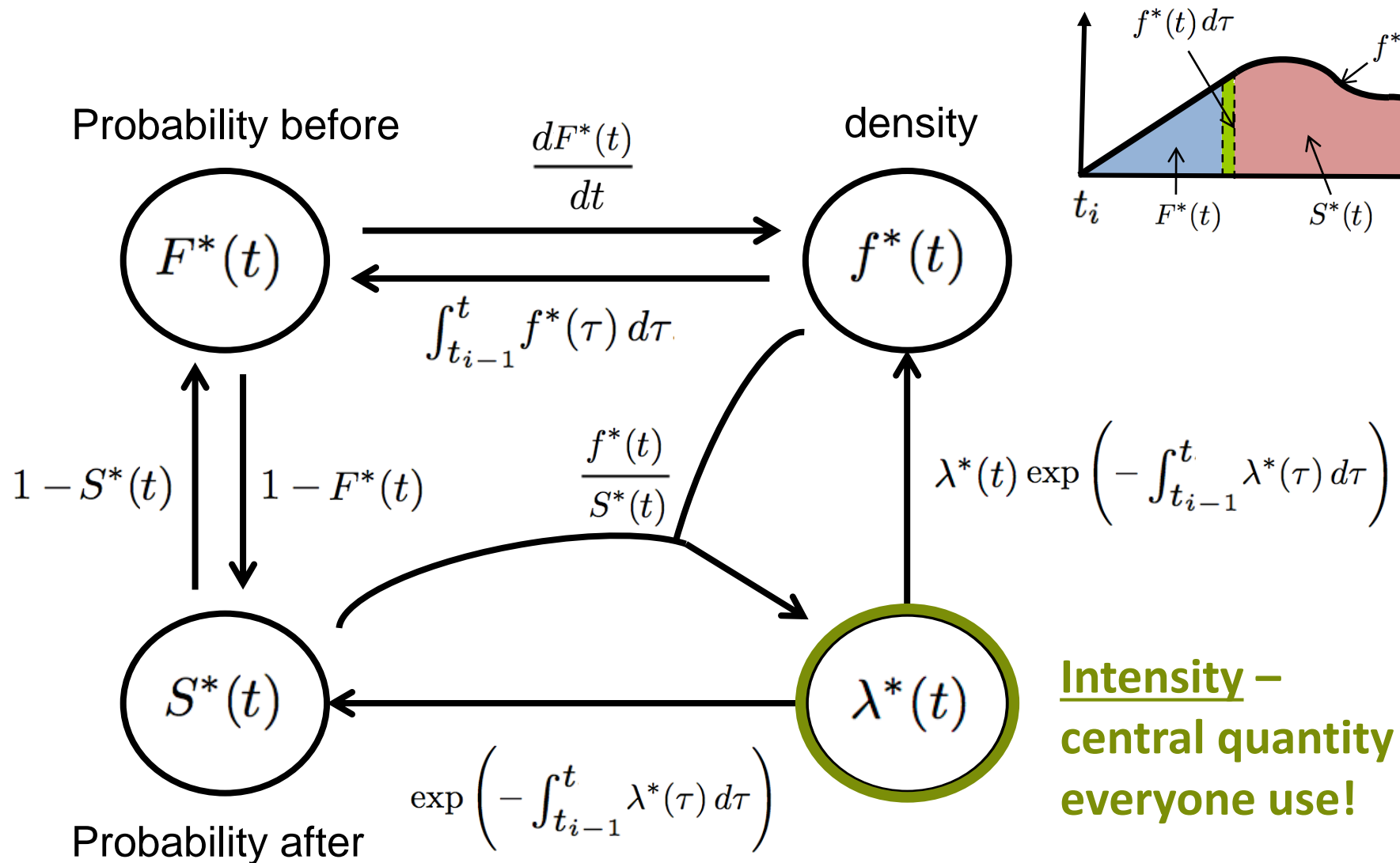
Sampling using inversion sampling:

$$t \sim \underbrace{\mu \exp(-\mu(t - t_3))}_{f_t^*(t)} \quad \Rightarrow \quad t = \underbrace{-\frac{1}{\mu} \log(1 - u)}_{F_t^{-1}(u)} + t_3$$

$\text{Uniform}(0, 1)$   
 $\downarrow$   
 $u$



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



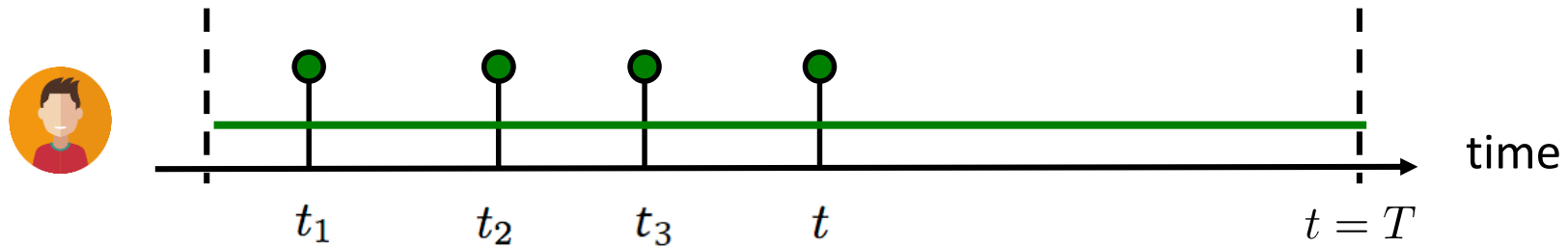


# **Representation: Temporal Point Processes**

## **Examples of processes**

# Poisson process

Coarse approximation of many real-life processes



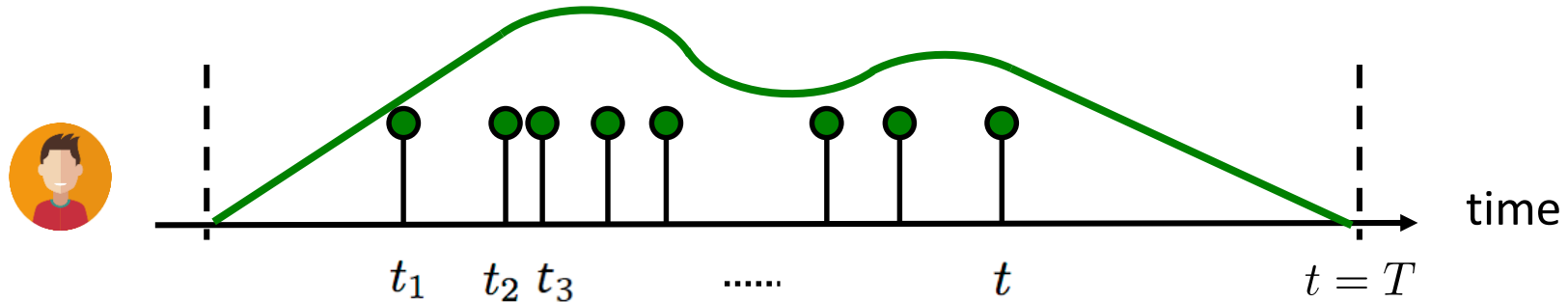
Intensity of a Poisson process

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

Observations:

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

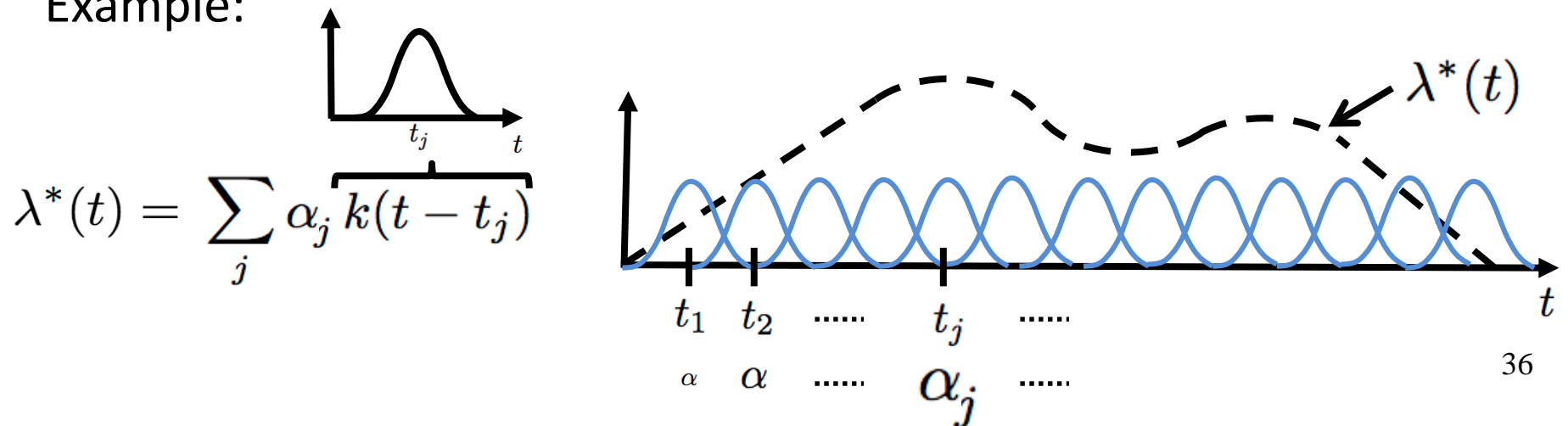
# Inhomogeneous Poisson process



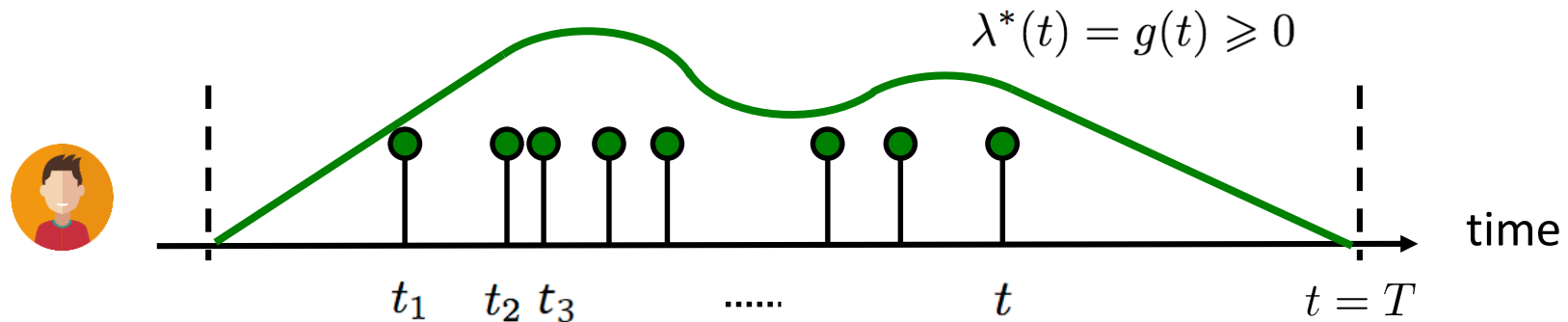
Intensity of an inhomogeneous Poisson process

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

Example:



# Fitting from inhomogeneous Poisson



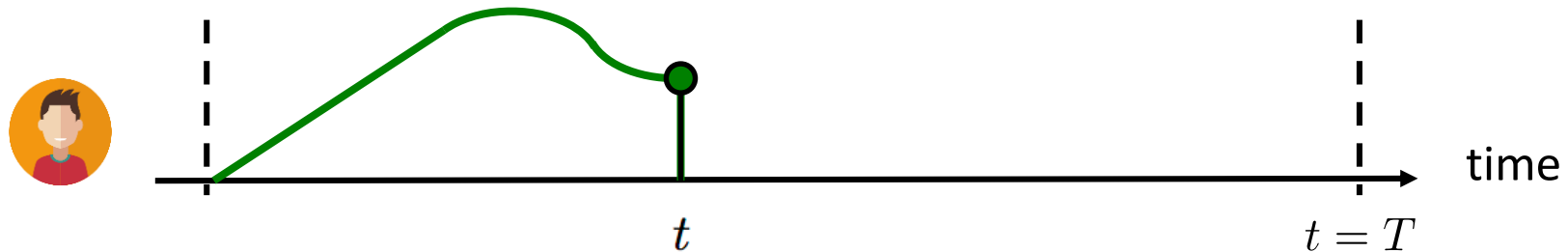
Fitting by maximum log-likelihood:

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

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



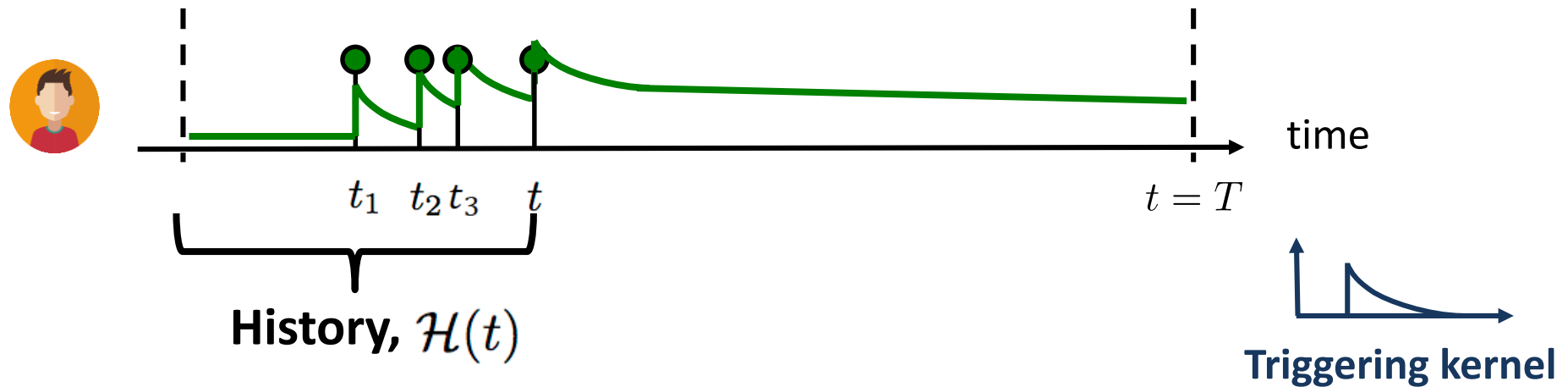
Intensity of a terminating (or survival) process

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

Observations:

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

# Self-exciting Hawkes process



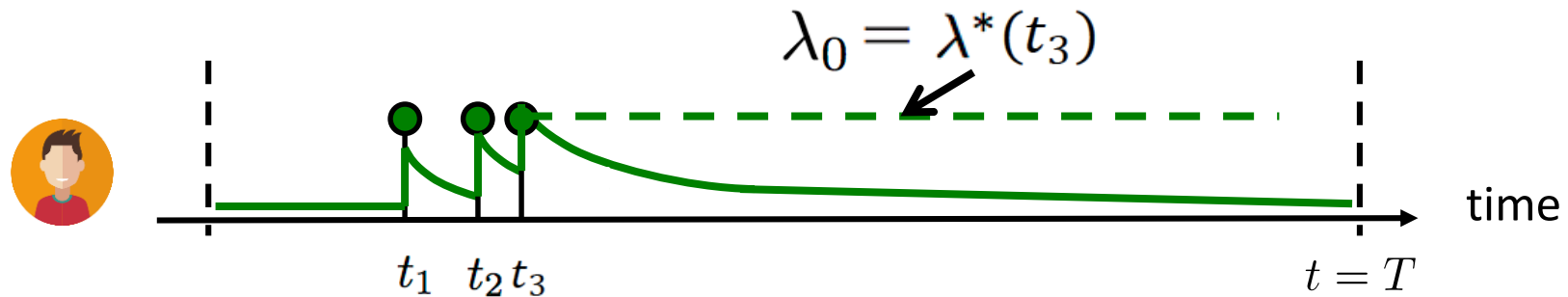
Intensity of self-exciting  
(or Hawkes) process:

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

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:

$$\text{maximize}_{\mu, \alpha} \left\{ \sum_{i=1}^n \log \lambda^*(t_i) - \int_0^T \lambda^*(\tau) d\tau \right\}$$

The max. likelihood is jointly convex in  $\mu$  and  $\alpha$

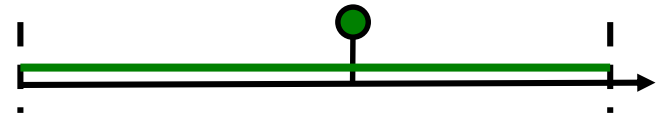


# Summary

## Building blocks to represent different dynamic processes:

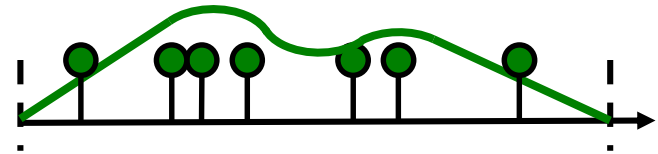
Poisson processes:

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



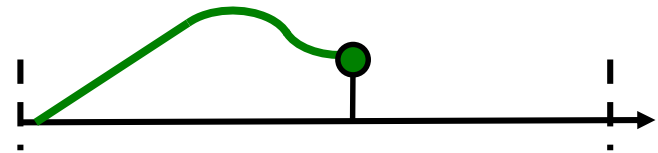
Inhomogeneous Poisson processes:

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



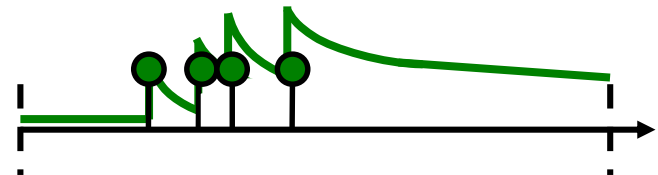
Terminating point processes:

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



Self-exciting point processes:

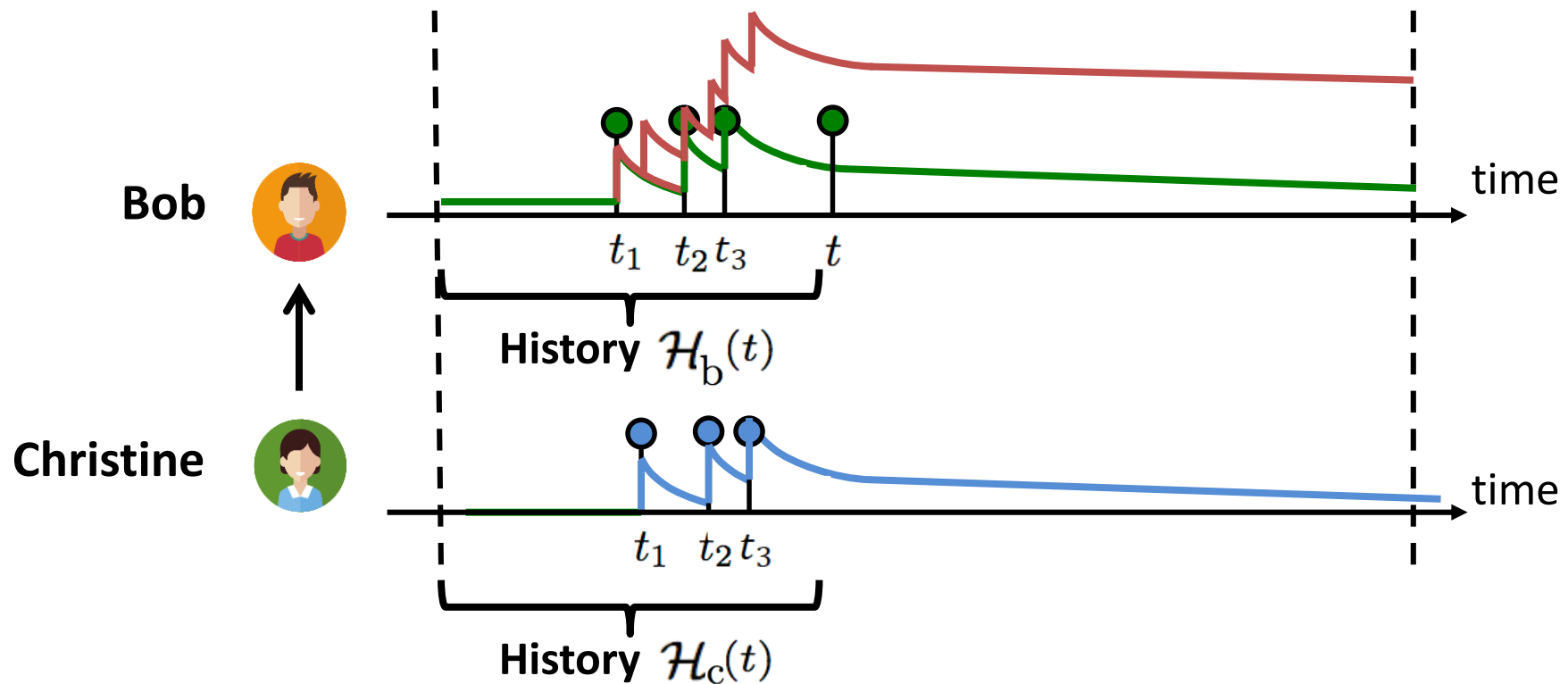
$$\lambda^*(t) = \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i)$$





# **Temporal Point Processes: other ideas**

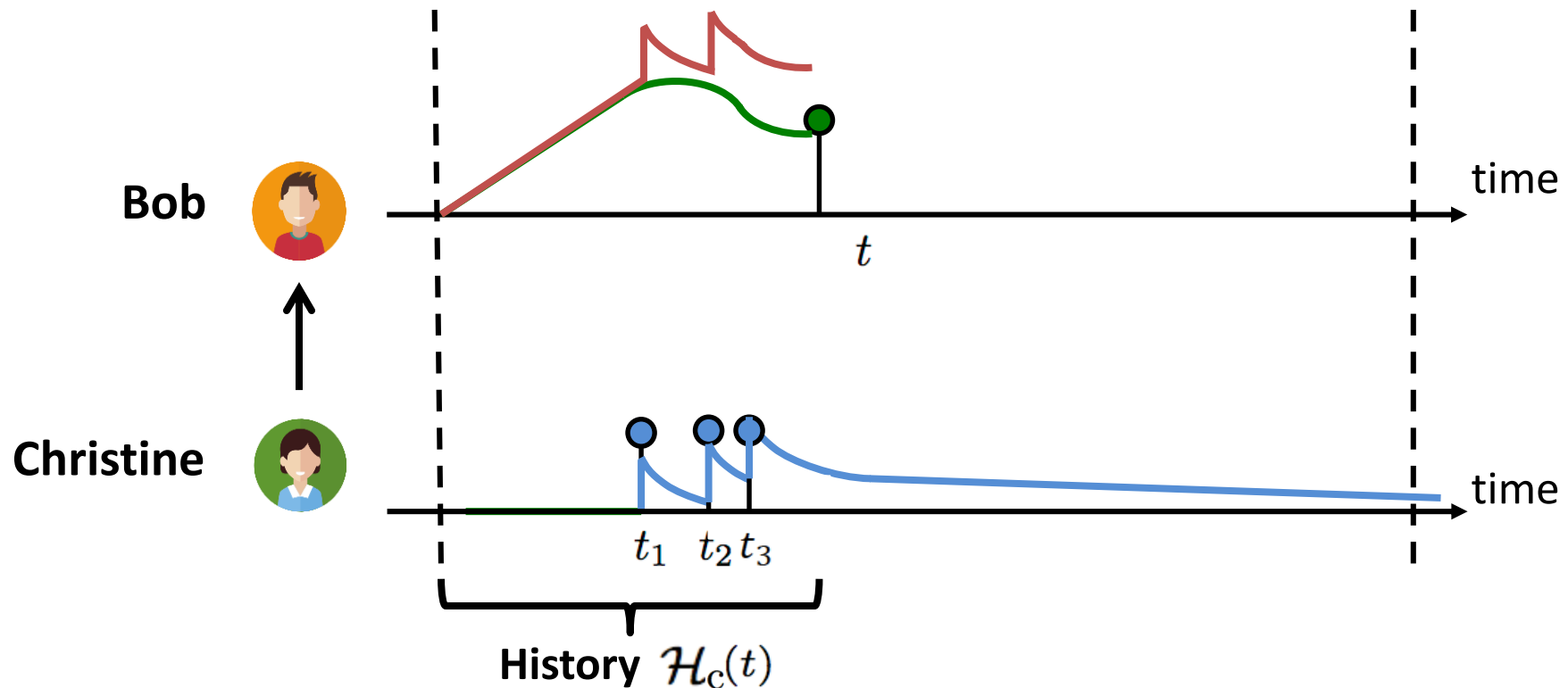
# Mutually exciting process



Clustered occurrence affected by neighbors

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

# Mutually exciting terminating process



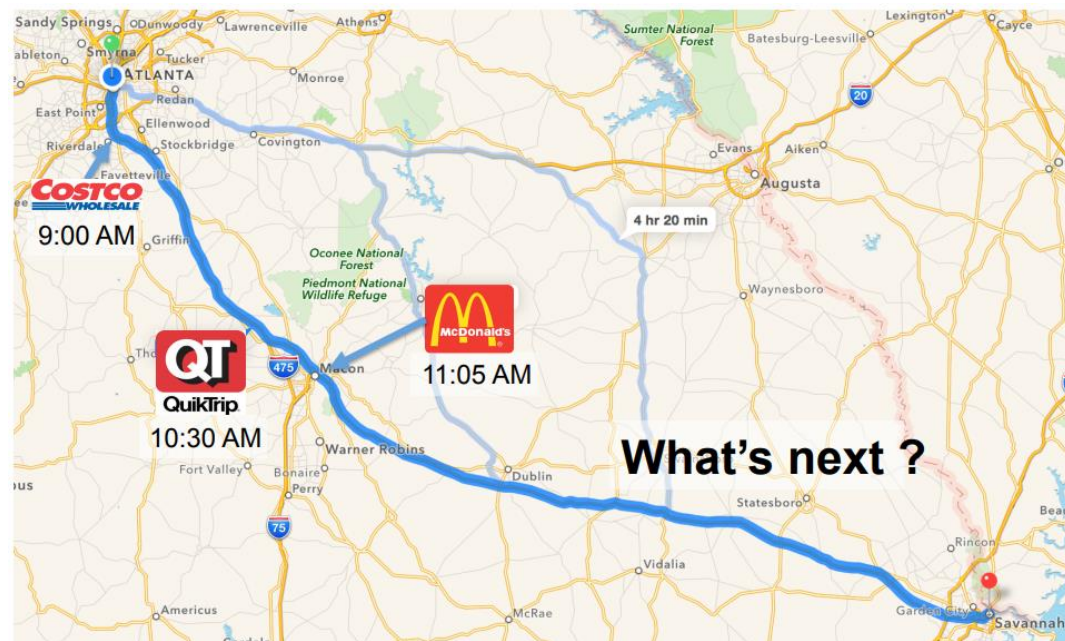
Clustered occurrence affected by neighbors

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

# Marked temporal point processes

Marked temporal point process:

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

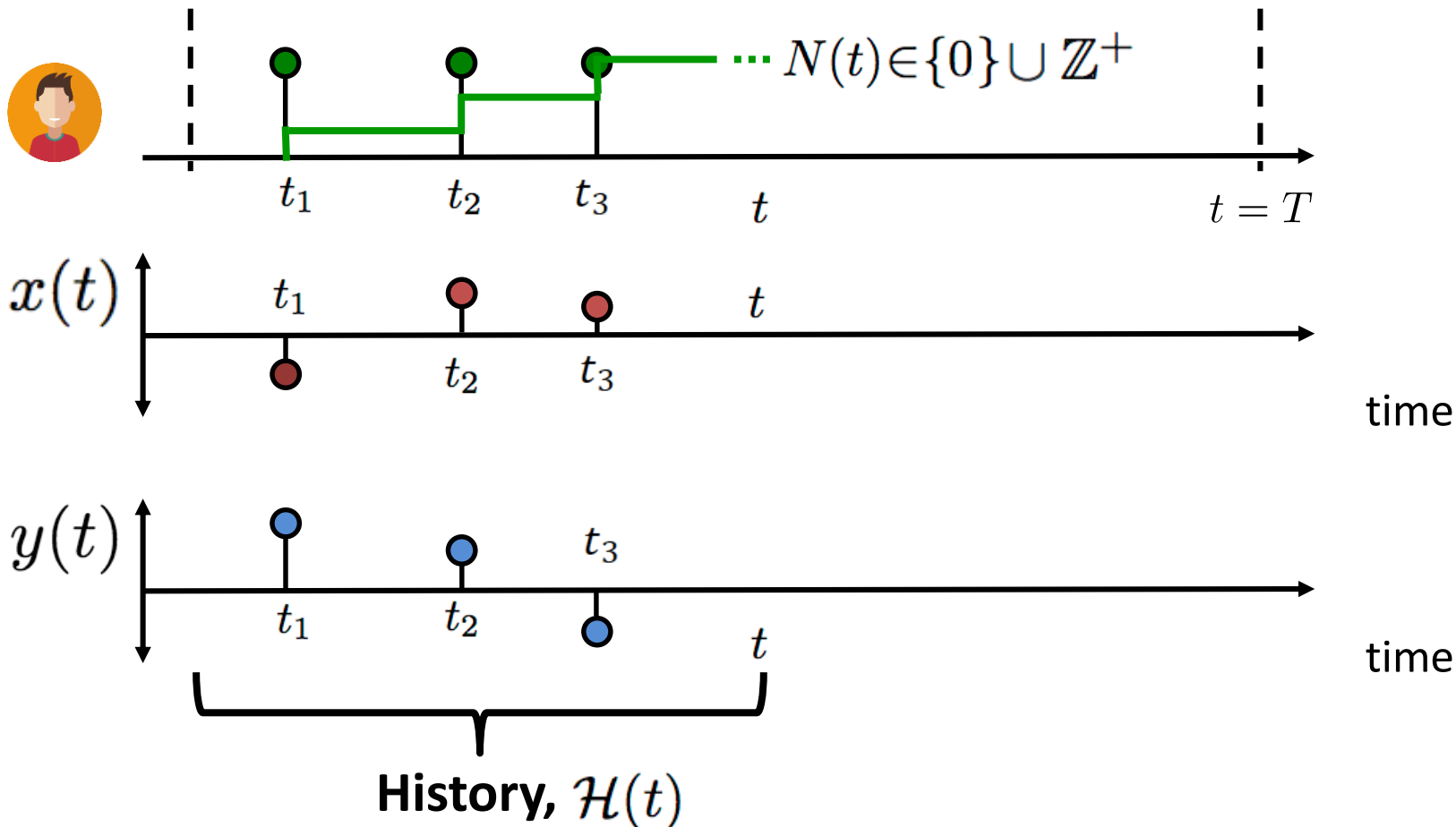


Given the trace of past locations and time, can we predict the location and time of the next stop?

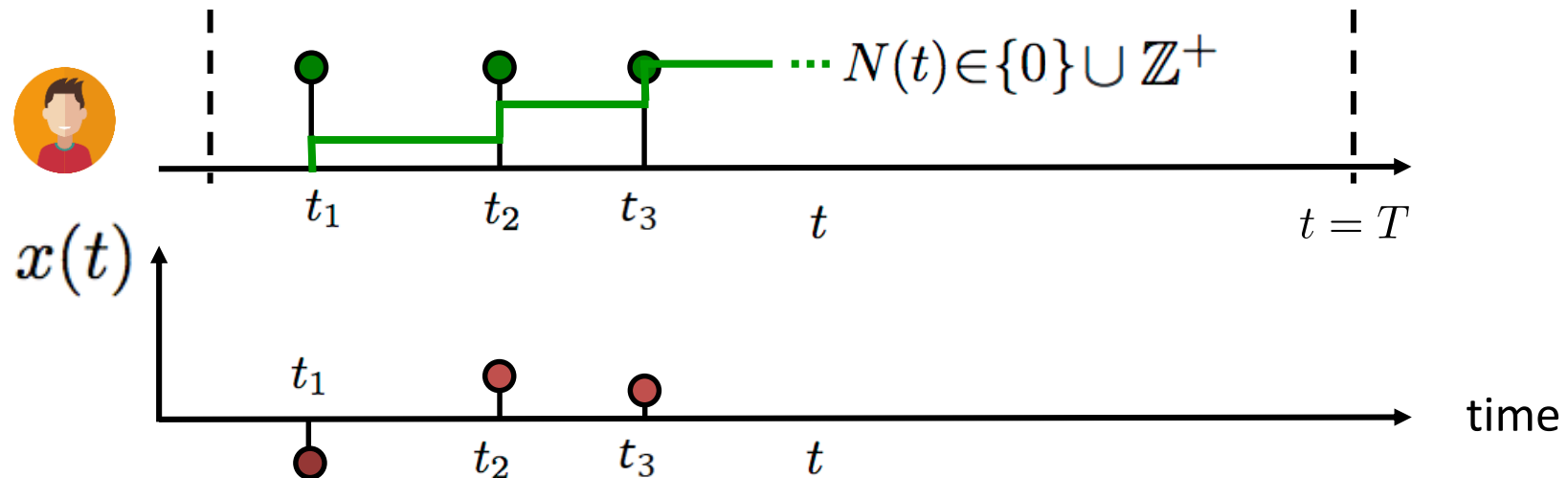
# Marked temporal point processes

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.)



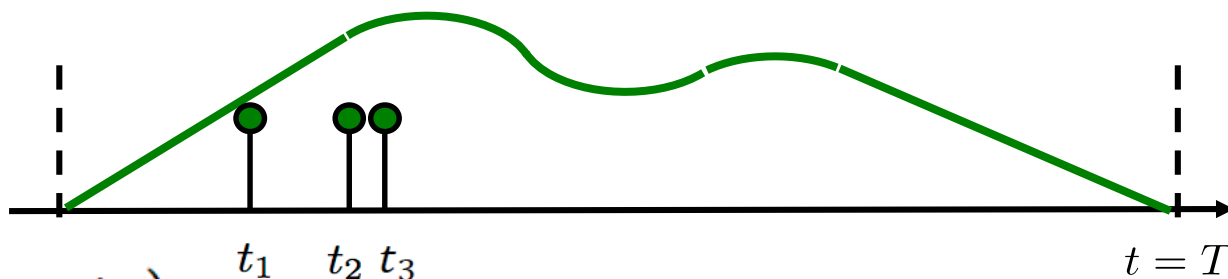
# **Models & Inference: Neural networks for the win**



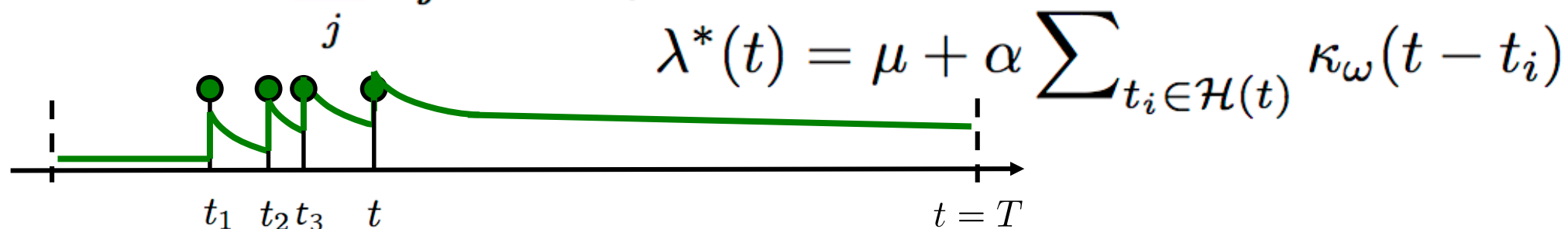
# Neural networks for the win

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

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



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

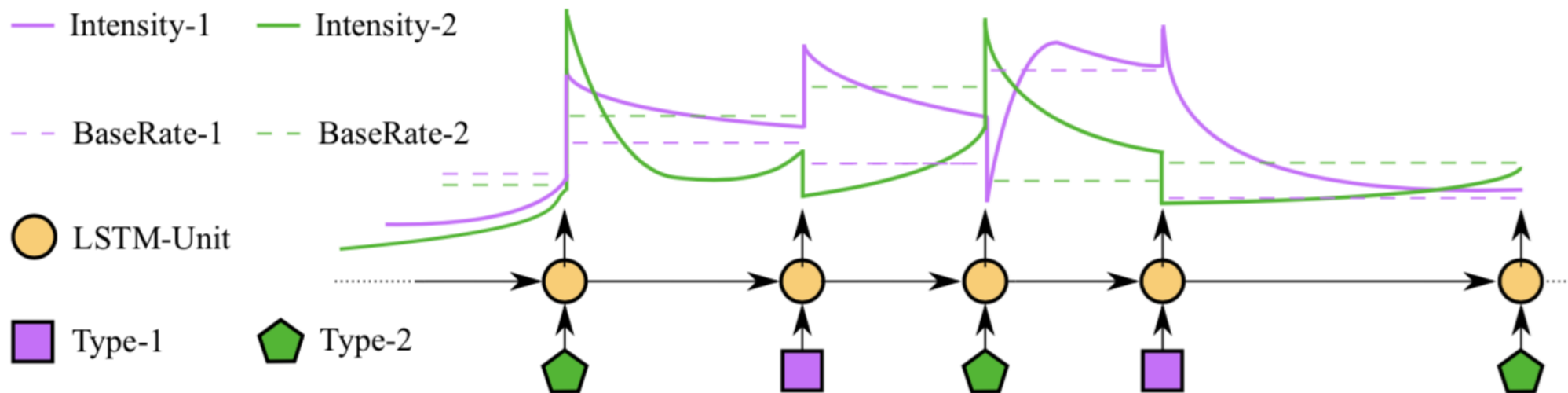


$$\lambda^*(t) = \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i)$$

Recent works make use of **RNNs** to capture more complex dynamics

# Neural Hawkes process

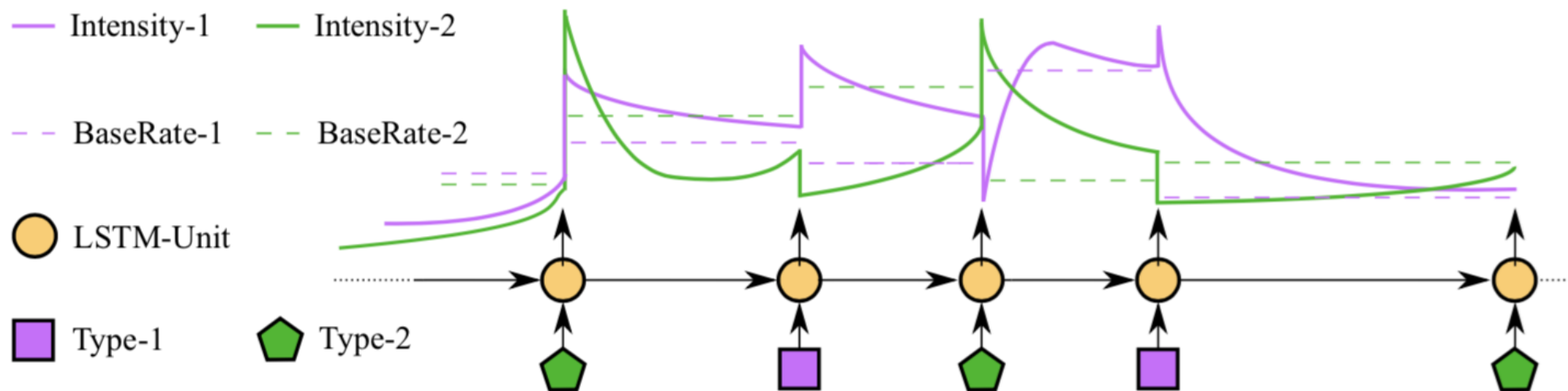
- 1) History effect does not need to be additive
- 2) Allows for complex memory effects (such as delays)



# Neural Hawkes process: RNNs

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

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



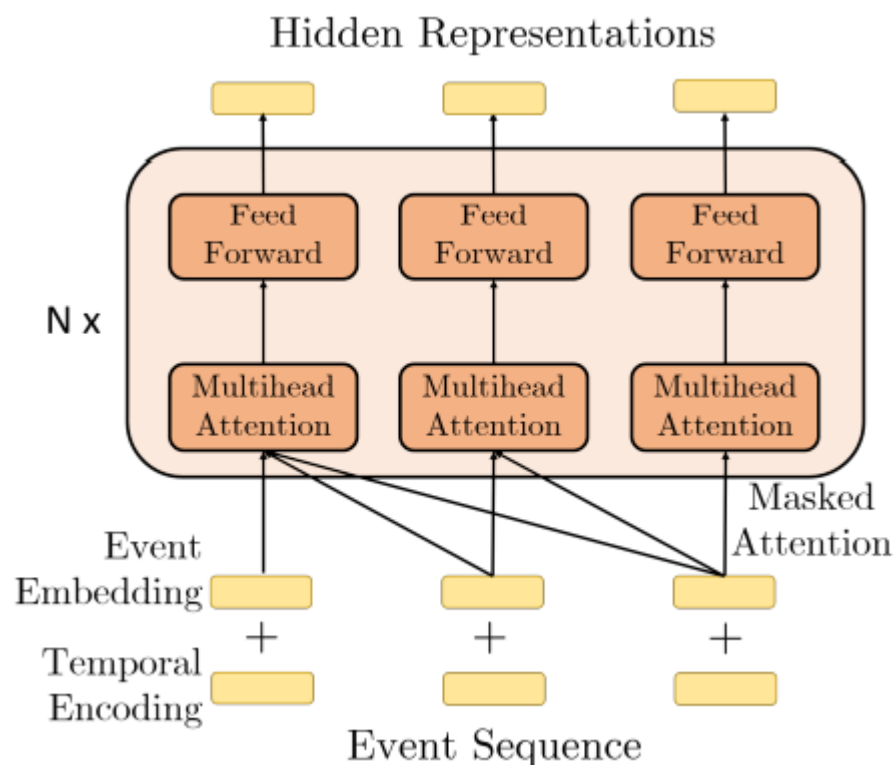
# Neural Hawkes process: Transformers

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

$$p(t|\mathcal{H}_t) = \lambda(t|\mathcal{H}_t) \exp \left( - \int_{t_j}^t \lambda(\tau|\mathcal{H}_\tau) d\tau \right),$$

$$\hat{t}_{j+1} = \int_{t_j}^{\infty} t \cdot p(t|\mathcal{H}_t) dt,$$

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





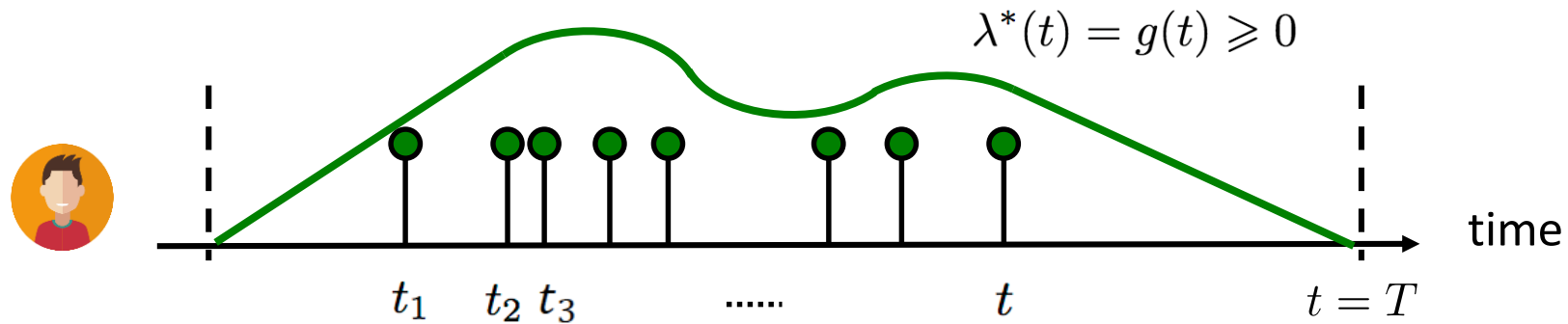
# References

# References

- Lifestream library by Sber <https://github.com/dlillb/pytorch-lifestream>
- Github with exercises  
<https://colab.research.google.com/drive/1Wc6aUNUZpE64egk4XwPT54v3sWDS495u?usp=sharing>
- Competition to predict gender based on transactions  
<https://www.kaggle.com/competitions/transactions/overview>  
<https://colab.research.google.com/drive/1tAqT4B5H9mkCiA3AilqUdIPOLpjEsXIK?usp=sharing>
- ICML tutorial on temporal point processes <https://learning.mpi-sws.org/tpp-icml18/>
- Lecture Notes: Temporal Point Processes and the Conditional Intensity Function by J. Rasmussen <https://arxiv.org/pdf/1806.00221.pdf>
- Zuo, Simiao, et al. "Transformer Hawkes process." *ICML*, 2020.

# Bonus II: Sampling

# Fitting & sampling from inhomogeneous Poisson



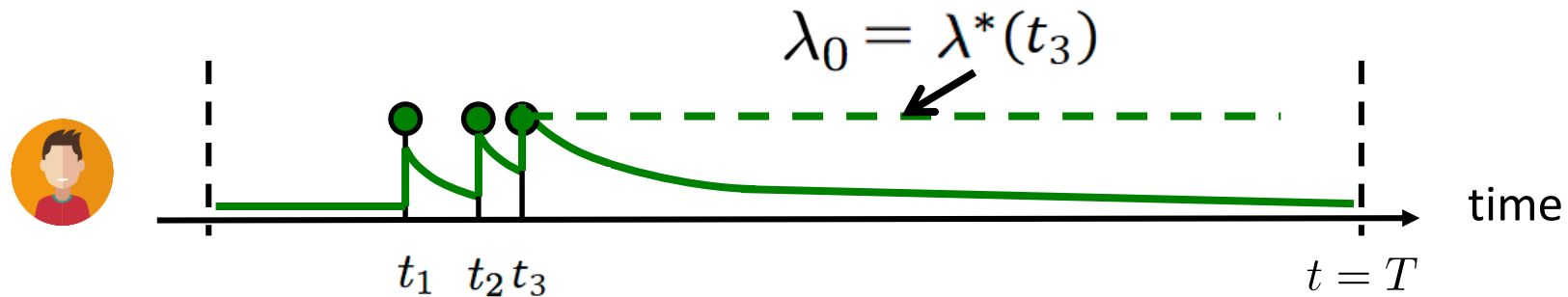
Fitting by maximum likelihood:  $\underset{g(t)}{\text{maximize}} \sum_{i=1}^n \log g(t_i) - \int_0^T g(\tau) d\tau$

Sampling using thinning (reject. sampling) + inverse sampling\*:

1. Sample  $t$  from Poisson process with intensity  $\mu$  using inverse sampling
  2. Generate  $u_2 \sim \text{Uniform}(0, 1)$
  3. Keep the sample if  $u_2 \leq g(t) / \mu$
- } Keep sample with prob.  $g(t) / \mu$



# Fitting a Hawkes process from a recorded timeline



Fitting by maximum likelihood:

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

Sampling using thinning (reject. sampling) + inverse sampling\*:

Key idea: the maximum of the intensity  $\lambda_0$  changes over time