

Learning the Joint Representation of Heterogeneous Temporal Events for Clinical Endpoint Prediction

Ming Zhang

School of Electronics Engineering and Computer Science

mzhang@net.pku.edu.cn

<http://net.pku.edu.cn/dlib/mzhang>



北京大學

<http://net.pku.edu.cn/dlib/>



<http://csrankings.org/>

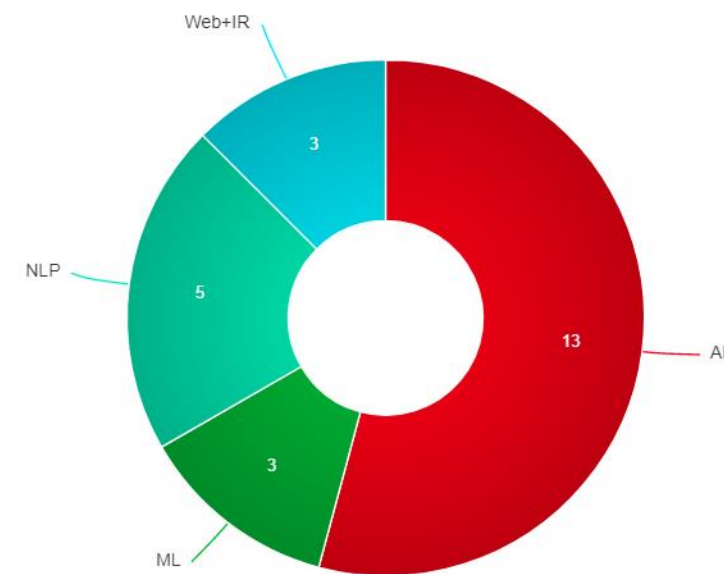
Ming Zhang 0004 AI

24

5.2

Ming Zhang 0004

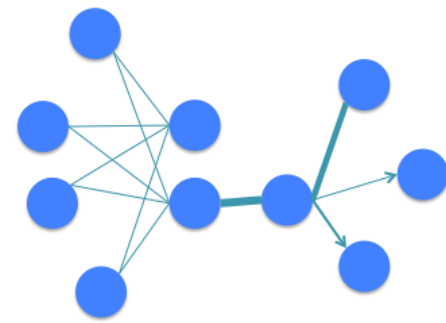
Publication Profile



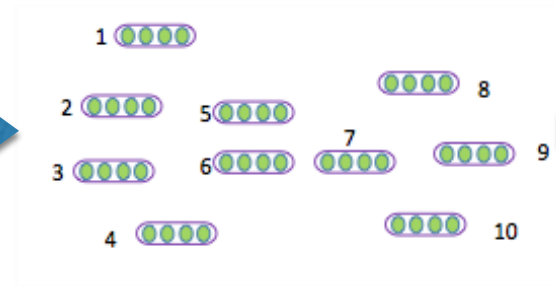
Ongoing Projects

- Co-PI, Multi-source and Heterogeneous Big Data Fusion Method Research Based on Whole Process Smart Health Management Decision, NSFC Key project, 2017-2020.
- PI, Machine Learning Models Based on Knowledge Graph and Deep Neural Network, Beijing Science and Technology Commission, 2018-2019.
- PI, A Knowledge Graphs Assisted General Framework to Construct Automatic Human-Computer Dialogue Systems for Vertical Domains, NSFC, 2018-2021
- PI, Research on user retention in massive open online courses, NSFC, 2015-2018

Learning Representations of Large-Scale Networks



Network

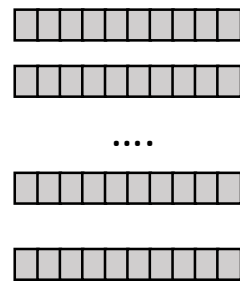


Node representations

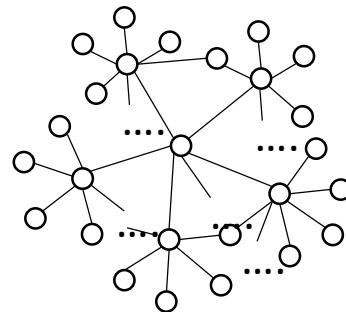


- Node Classification
- Node Clustering
- Link Prediction
- Recommendation
- ...

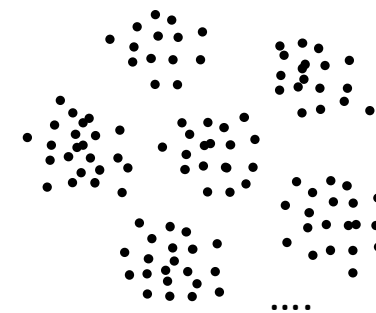
- E.g., Facebook social network -> user representations (features)-> friend recommendation



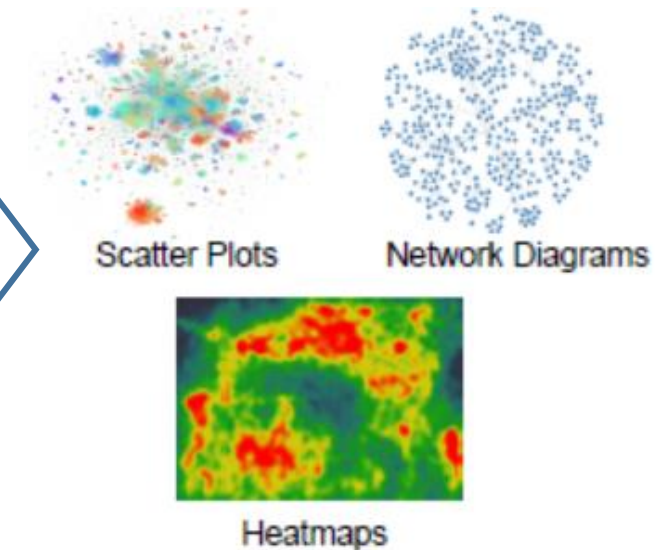
High-dimensional Data



Networks



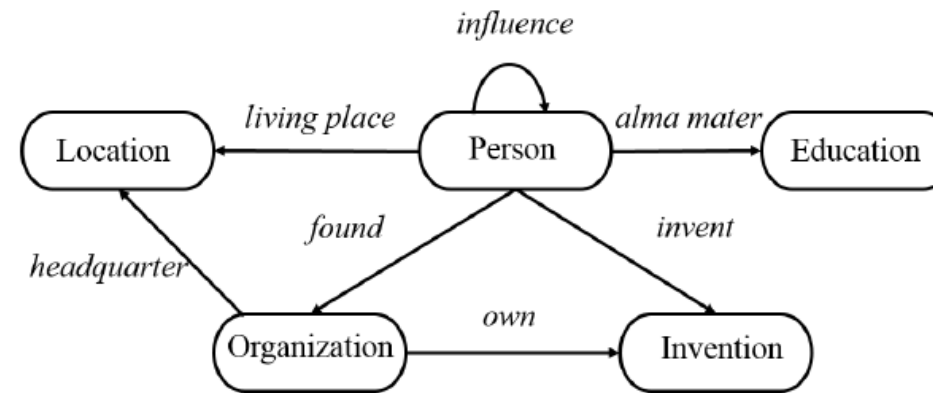
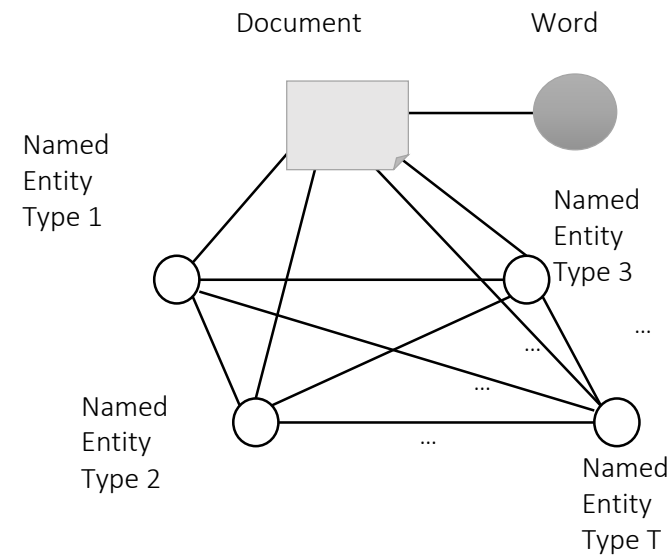
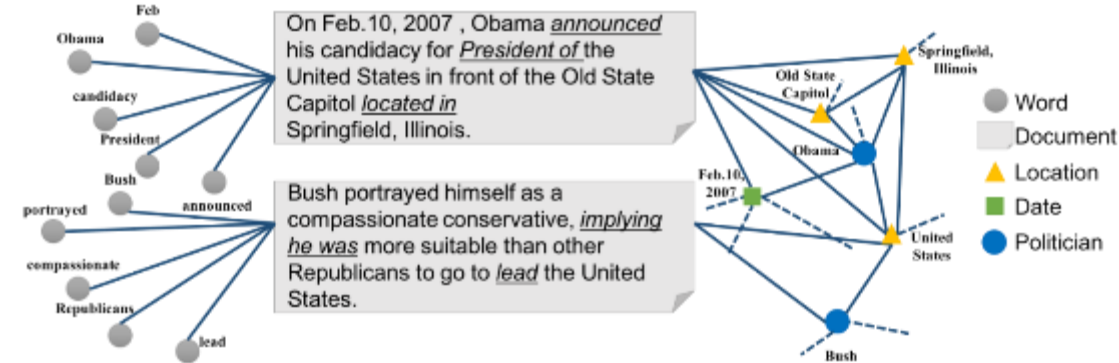
2D/3D Layout



Network Embedding

- Jian Tang, Meng Qu, Mingzhe Wang, **Ming Zhang**, Jun Yan, Qiaozhu Mei. LINE: Large-scale Information Network Embedding. WWW 2015, 1067-1077. **citations 747**
- Jian Tang, Zhaoshi Meng, XuanLong Nguyen, Qiaozhu Mei, **Ming Zhang**, Understanding the Limiting Factors of Topic Modeling via Posterior Contraction Analysis, The 31st International Conference on Machine Learning (**ICML2014**), **Best paper**, 2014.6.21-2014.6.26, **citations 123**
- Jian Tang, Jingzhou Liu, **Ming Zhang**, Qiaozhu Mei, Visualizing Large-scale and High-dimensional Data, **WWW 2016**. 04.11-2016.04.15, **Best paper runner-up**, **citations 92**
- [Jian Tang](#), Ming Zhang, [Qiaozhu Mei](#): One theme in all views: modeling consensus topics in multiple contexts. [KDD 2013](#): 5-13, **citations 39**
- [Meng Qu](#), [Jian Tang](#), [Jingbo Shang](#), [Xiang Ren](#), Ming Zhang, [Jiawei Han](#): An Attention-based Collaboration Framework for Multi-View Network Representation Learning. [CIKM 2017](#): 1767-1776, **citations 8**

World Knowledge Representation: Heterogeneous Information Network (HIN)



Incorporating World Knowledge to Heterogeneous Information Networks

- [1] Chenguang Wang, Yizhou Sun, Yanglei Song, Jiawei Han, Yangqiu Song, Lidan Wang, and **Ming Zhang**: RelSim: RelSim: Relation Similarity Search in Schema-Rich Heterogeneous Information Networks. Proc. 2016 SIAM Int. Conf. on Data Mining (**SDM'16**).**citations 11**
- [2] Chenguang Wang, Yangqiu Song, Haoran Li, **Ming Zhang**, and Jiawei Han: Text Classification with Heterogeneous Information Network Kernels. Proc. 2016 AAAI Conf. on Artificial Intelligence (**AAAI'16**).**他34**
- [3] Chenguang Wang, Yangqiu Song, Haoran Li, **Ming Zhang**, and Jiawei Han: KnowSim: A Document Similarity Measure on Structured Heterogeneous Information Networks. Proc. of 2014 IEEE Int. Conf. on Data Mining (**ICDM'15**).**citations 32**
- [4] Chenguang Wang, Yangqiu Song, Ahmed El-Kishky, Dan Roth, **Ming Zhang**, and Jiawei Han: Incorporating World Knowledge to Document Clustering via Heterogeneous Information Networks. Proc. 2015 ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (**KDD'15**).**citations 30**
- [5] Chenguang Wang, Yangqiu Song, Dan Roth, Chi Wang, Jiawei Han, Heng Ji, and **Ming Zhang**: Constrained Information-Theoretic Tripartite Graph Clustering to Identify Semantically Similar Relations. Proc. 2015 Int. Joint Conf. on Artificial Intelligence (**IJCAI'15**).**citations 13**
- [6] Yangqiu Song, Chenguang Wang, Ming Zhang and Hailong Sun: Spectral Label Refinement for Noisy and Missing Text Labels. Proc. 2015 AAAI Conf. on Artificial Intelligence (**AAAI'15**).**citations 6**
- [7] [Chenguang Wang](#), [Yangqiu Song](#), [Dan Roth](#), Ming Zhang, [Jiawei Han](#): World Knowledge as Indirect Supervision for Document Clustering. [TKDD 11\(2\)](#): 13:1-13:36 (2016).**citations 4**
- [9] [He Jiang](#), [Yangqiu Song](#), [Chenguang Wang](#), Ming Zhang, [Yizhou Sun](#): Semi-supervised Learning over Heterogeneous Information Networks by Ensemble of Meta-graph Guided Random Walks. [IJCAI 2017](#): 1944-1950.**citations 6**
- [8] [Chenguang Wang](#), [Yangqiu Song](#), [Haoran Li](#), [Yizhou Sun](#), Ming Zhang, [Jiawei Han](#): Distant Meta-Path Similarities for Text-Based Heterogeneous Information Networks. [CIKM2017](#): 1629-1638.**citations 6**

Learning the Joint Representation of Heterogeneous Temporal Events for EHR

- LuchenLiu, Jianhao Shen, Ming Zhang, Zichang Wang, Jian Tang: Learning the Joint Representation of Heterogeneous Temporal Events for Clinical Endpoint Prediction. AAAI 2018
- LuchenLiu, Haoran Li, Jianhao Shen, Ming Zhang, Zichang Wang, Jian Tang: Modeling Temporal Events with Heterogeneous Attributes by Extracting Latent Groups for Clinical Outcome Prediction. AAAI 2019 in submission
- LuchenLiu, Jianhao Shen, Ming Zhang, Zichang Wang, Jian Tang: Learning Hierarchical Representations of Heterogeneous Event Sequences. JAMIA in submission

Outline

- Background—endpoint prediction in EHR
- Heterogeneous temporal events
- Model
- Experiments

Outline

- Background—endpoint prediction in EHR
- Heterogeneous temporal events
- Model
- Experiments

Back



图片来源：视觉中国 www.vcg.com



中国人民解放军总医院海南分院
检验报告单

姓名: 刘薇 女 18岁 费别: 全费 门诊号: HN301607
住院号: 014835

(K) 血清钾等 申请序号: 1414259735 报告时间: 2014-02-26 09:33
检验科 检验人: 张自康 审核人: 邓心慧 工作单号: CH0006

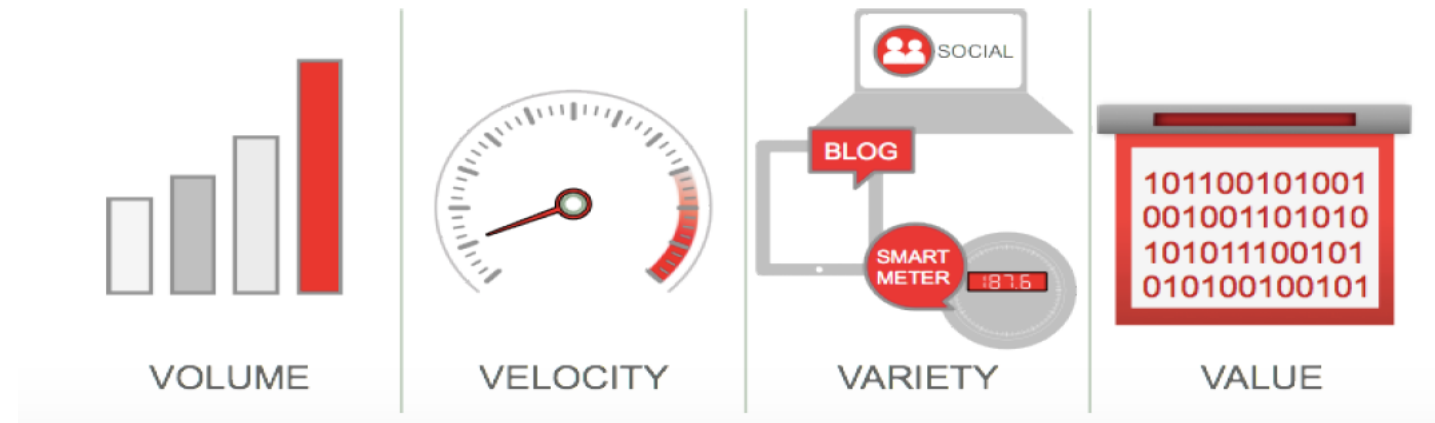
检验项目	结果	标志	参考值	单位
碱性磷酸酶	71		0~130	U/L
γ-谷氨酰转氨酶	57		0~60	U/L
葡萄糖	4.12		3.4~6.1	mmol/L
尿素	1.2		1.8~7.5	mmol/L
肌酐	30		30~110	umol/L
血清尿酸	260		104~444	umol/L
总胆红素	2.93		3.1~5.7	mmol/L
甘油三酯	0.95		0.4~1.7	mmol/L
载脂蛋白A1	1.13		1.0~1.6	g/L
载脂蛋白B	0.69		0.6~1.1	g/L
肌酸酐酶	13		2~250	U/L
乳酸脱氢酶	195		40~250	U/L
肌酸酐酶同工酶	8.5		0~24	U/L



HER (Electronic Health Record)



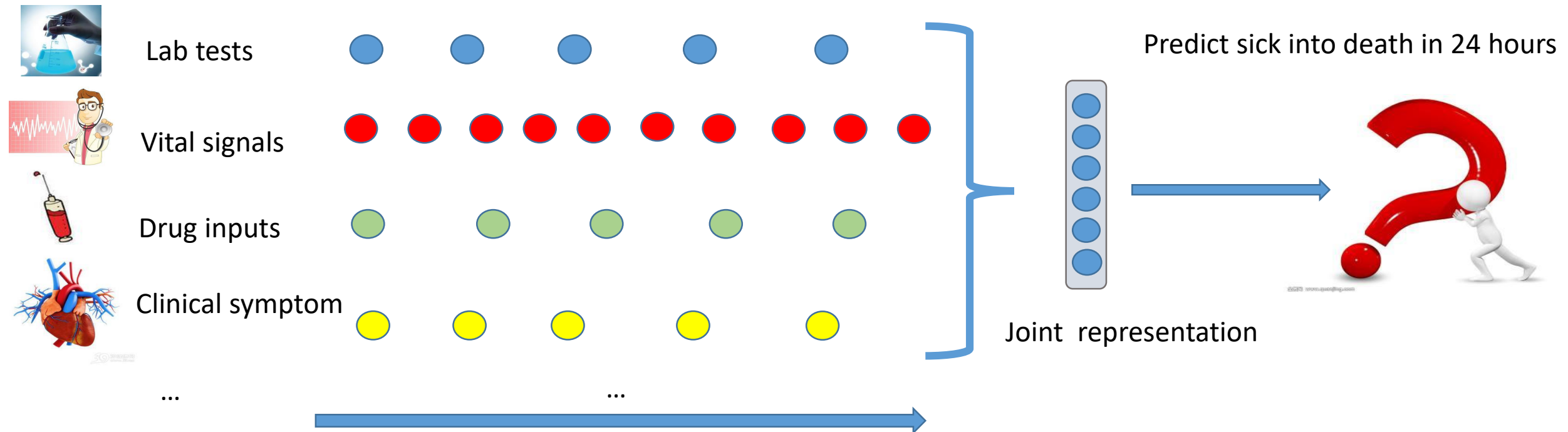
Background



- Challenges in today's healthcare system — — limited analysis ability for Big Data
 - Lab test results \neq disease mechanism
 - Statistical significance \neq clinical significance
 - A certain medicine can reduce the risk of heart attack by 34%
 - Real clinical risk are reduced by 1.4% (N Engl J med. 1987 Nov 12; 317(20):1237-45)
- The value of Medical Big Data
 - HER (**E**lectronic **H**ealth **R**ecord) ~ Clinical Endpoint (the target outcomes, e.g. death, symptom...)

Endpoint prediction based on EHR

- EHR events embedding
 - Clinical events → patient status representation reflect the disease mechanism
- clinical endpoint prediction
 - Clinical endpoint prediction → personalized diagnosis decision



Outline

- Background—endpoint prediction in EHR
- **Heterogeneous temporal events**
- Model
- Experiments

Heterogeneous temporal events

Thousands of event types



Lab tests



Vital signals



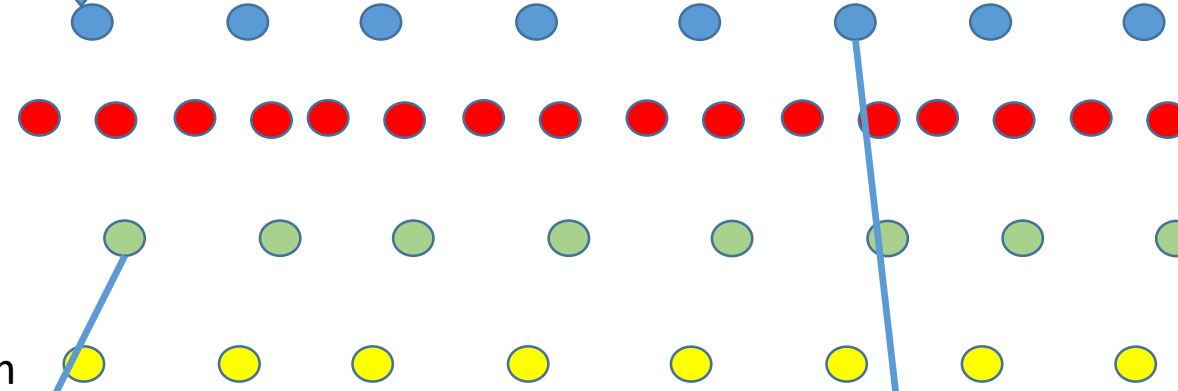
Drug inputs



Clinical symptom

...

event



...

Drug input



Solution type

dosage

rate

Lab test



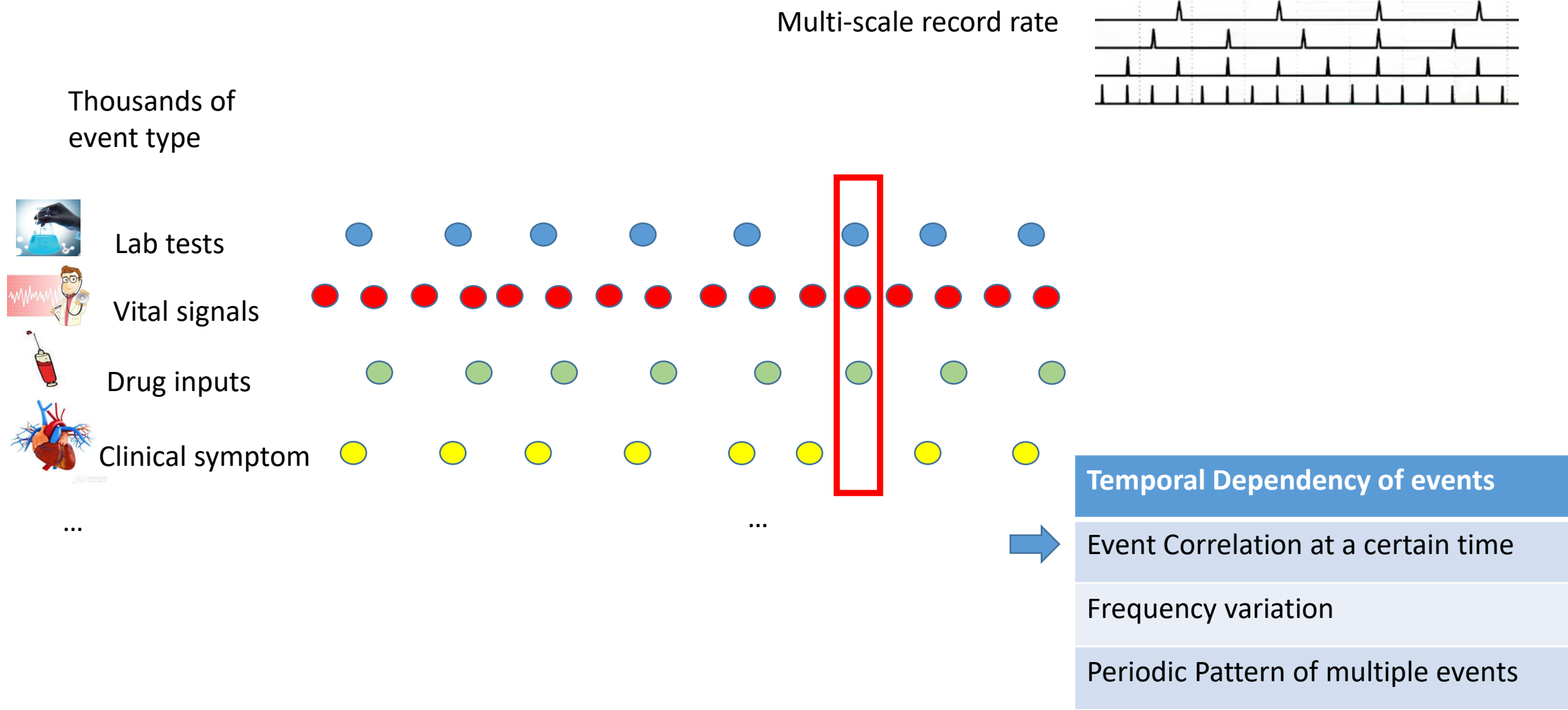
Abnormal
or normal

Result
index

Event ID	Event Type	Time	Value	Unit	Category
1	Lab test	2023-01-01 10:00	120	mg/dL	Glucose
2	Vital signal	2023-01-01 10:05	98	%	SpO2
3	Drug input	2023-01-01 10:10	5	mg	Insulin
4	Clinical symptom	2023-01-01 10:15	1		Hypertension
5	Lab test	2023-01-01 10:20	115	mg/dL	Glucose
6	Vital signal	2023-01-01 10:25	95	%	SpO2
7	Drug input	2023-01-01 10:30	10	mg	Insulin
8	Clinical symptom	2023-01-01 10:35	2		Hypertension
9	Lab test	2023-01-01 10:40	110	mg/dL	Glucose
10	Vital signal	2023-01-01 10:45	92	%	SpO2
11	Drug input	2023-01-01 10:50	15	mg	Insulin
12	Clinical symptom	2023-01-01 10:55	3		Hypertension
13	Lab test	2023-01-01 11:00	105	mg/dL	Glucose
14	Vital signal	2023-01-01 11:05	90	%	SpO2
15	Drug input	2023-01-01 11:10	20	mg	Insulin
16	Clinical symptom	2023-01-01 11:15	4		Hypertension
17	Lab test	2023-01-01 11:20	100	mg/dL	Glucose
18	Vital signal	2023-01-01 11:25	88	%	SpO2
19	Drug input	2023-01-01 11:30	25	mg	Insulin
20	Clinical symptom	2023-01-01 11:35	5		Hypertension

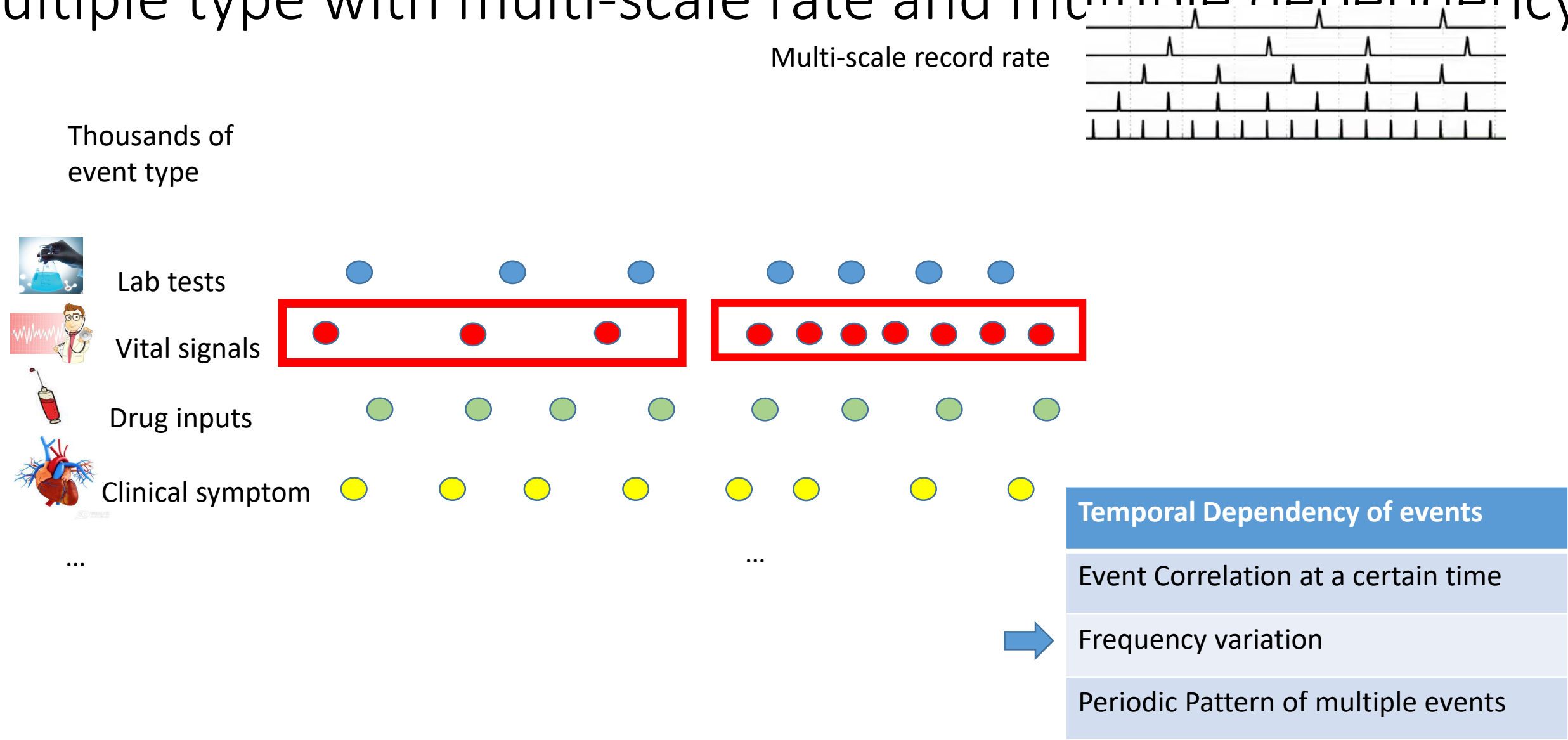
The frontier of heterogeneous temporal events

----multiple type with multi-scale rate and multiple dependency



The frontier of heterogeneous temporal events

----multiple type with multi-scale rate and multiple dependency



The frontier of heterogeneous temporal events

----multiple type with multi-scale rate and multiple dependency

Thousands of
event type



Lab tests



Vital signals

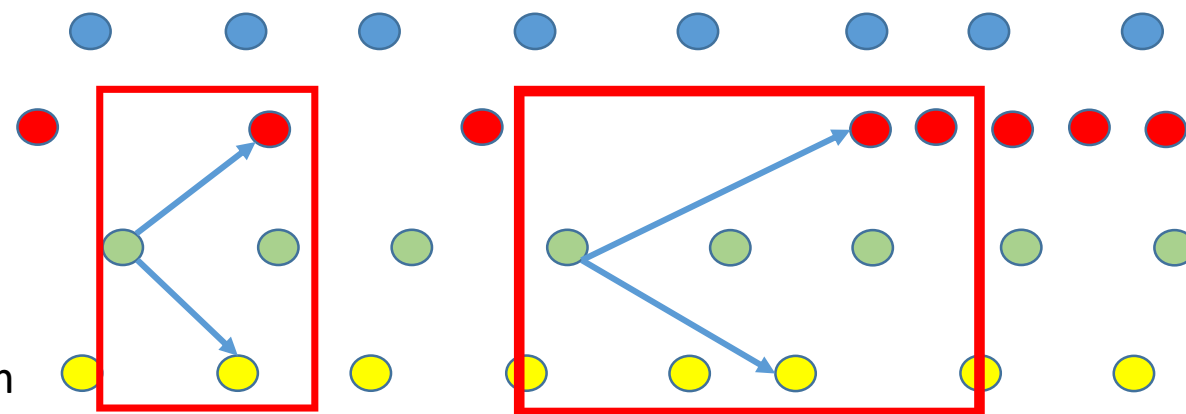


Drug inputs

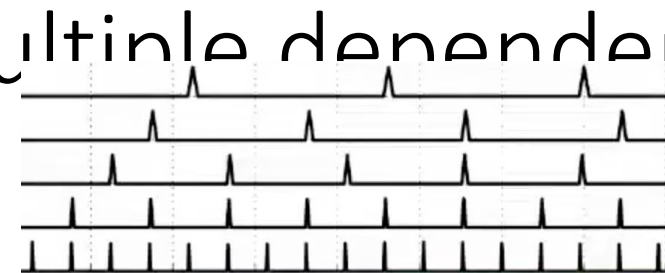


Clinical symptom

...



Multi-scale record rate



Temporal Dependency of events

Event Correlation at a certain time

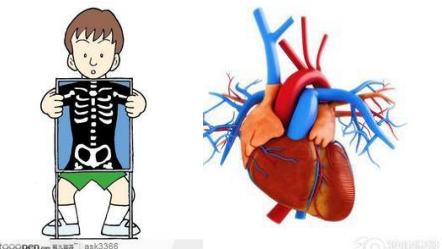
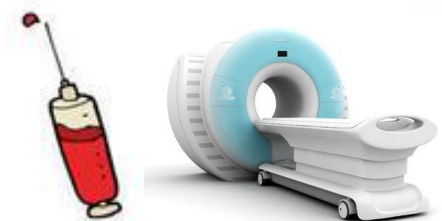
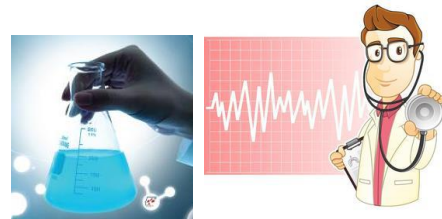
Frequency variation

Periodic Pattern of multiple events

Outline

- Background——endpoint prediction in EHR
- Heterogeneous temporal events
- **Model —— Heterogeneous Event LSTM**
- Experiments

Idea: work in cooperation
Asynchronously tracing important
information of related events



Input Stream



Neuron 1



"off"

Neuron 2



"off"

"off"

Neuron 3



"off"


Event gate

Asynchronously tracing important information of related events

- Event gate to decide whether or not to update hidden state

$$\tilde{c}_l = \mathbf{f}_l \circ \mathbf{c}_{l-1} + \mathbf{i}_l \circ \tanh(W_{cx}\mathbf{x}_l + W_{ch}\mathbf{h}_{l-1} + \mathbf{b}_c)$$
$$\mathbf{c}_l = \mathbf{j}_l \circ \tilde{c}_l + (1 - \mathbf{j}_l) \circ \mathbf{c}_{l-1}$$

- Event gate is controlled by the event type and time

$$\mathbf{j}_{s,t} = \mathbf{e}_s \circ \mathbf{k}_t$$
$$\mathbf{e}_s = \sigma(W_{em} \tanh(W_{ms}\mathbf{s} + \mathbf{b}_m) + \mathbf{b}_e)$$


The diagram shows a horizontal timeline with several vertical dashed lines representing time steps. A series of sharp upward spikes (events) are plotted on this timeline. Two blue arrows originate from the equation $\mathbf{j}_{s,t} = \mathbf{e}_s \circ \mathbf{k}_t$. One arrow points from \mathbf{e}_s to the first spike, and the other points from \mathbf{k}_t to the third spike, illustrating the element-wise multiplication of the event vector \mathbf{e}_s and the time vector \mathbf{k}_t to produce the event gate $\mathbf{j}_{s,t}$.

- Each neuron of the C vector refers to the status of a set of related **events** at a certain record **rate**

Heterogeneous Event LSTM(HE-LSTM)

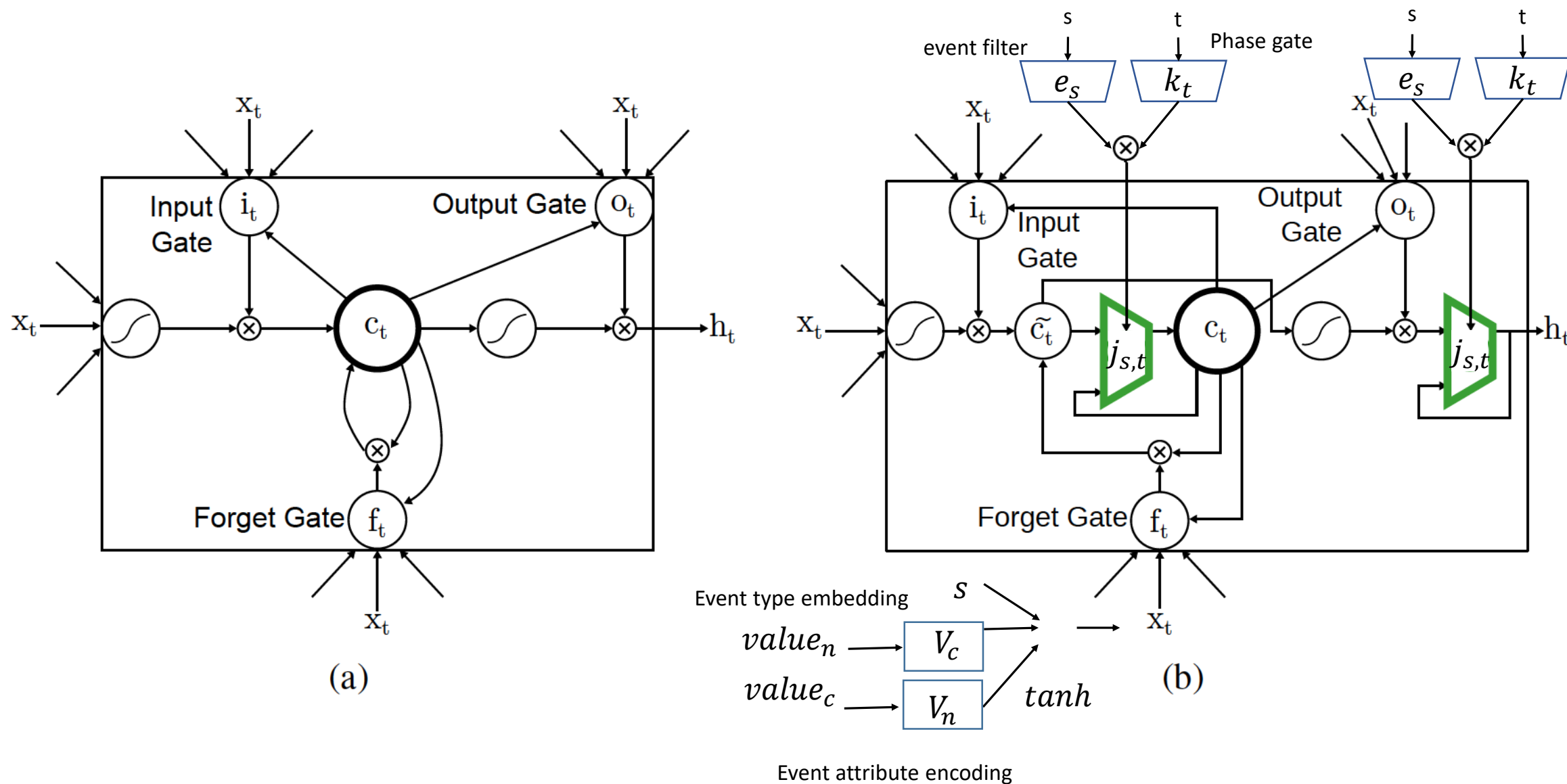


Illustration of the multiple dependency of related events

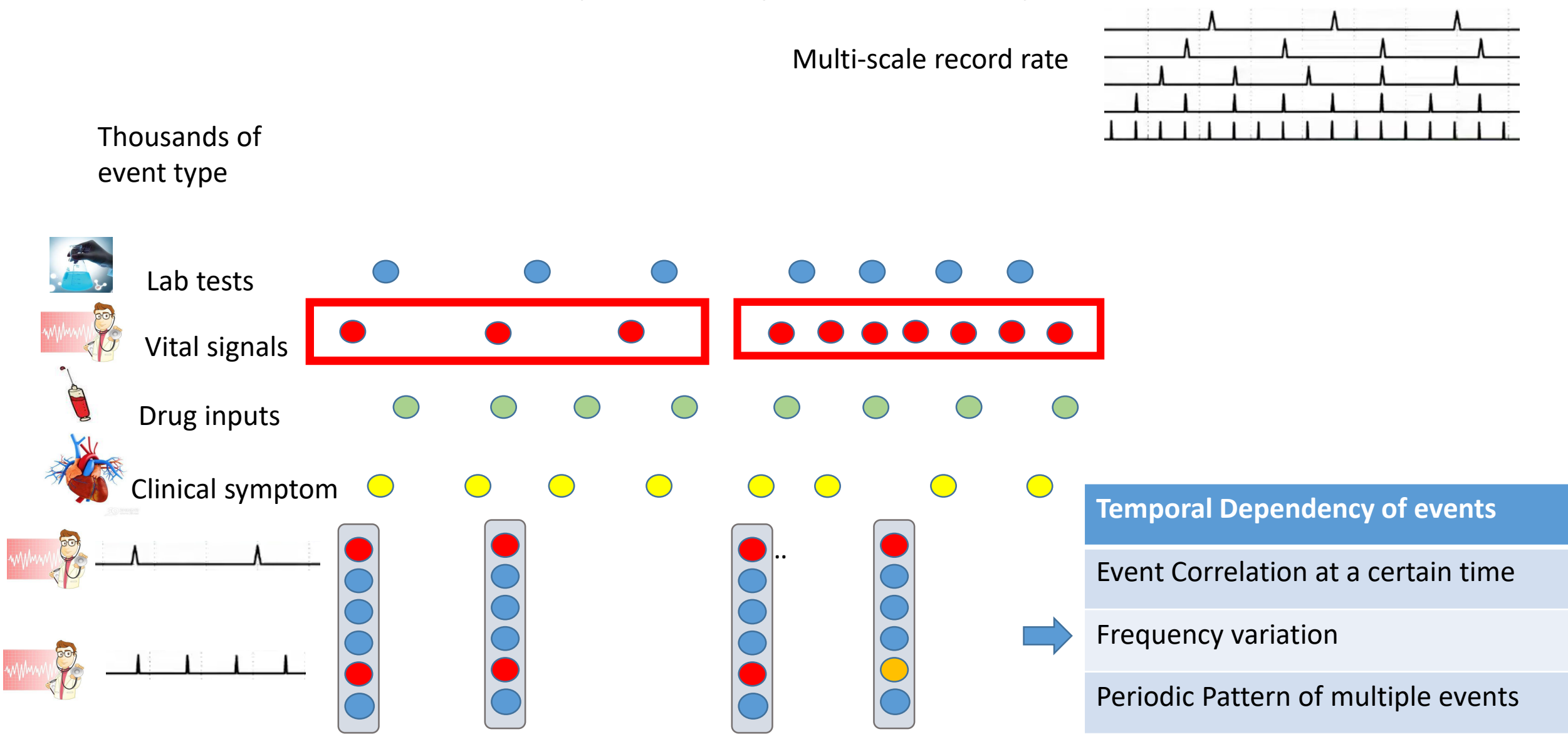
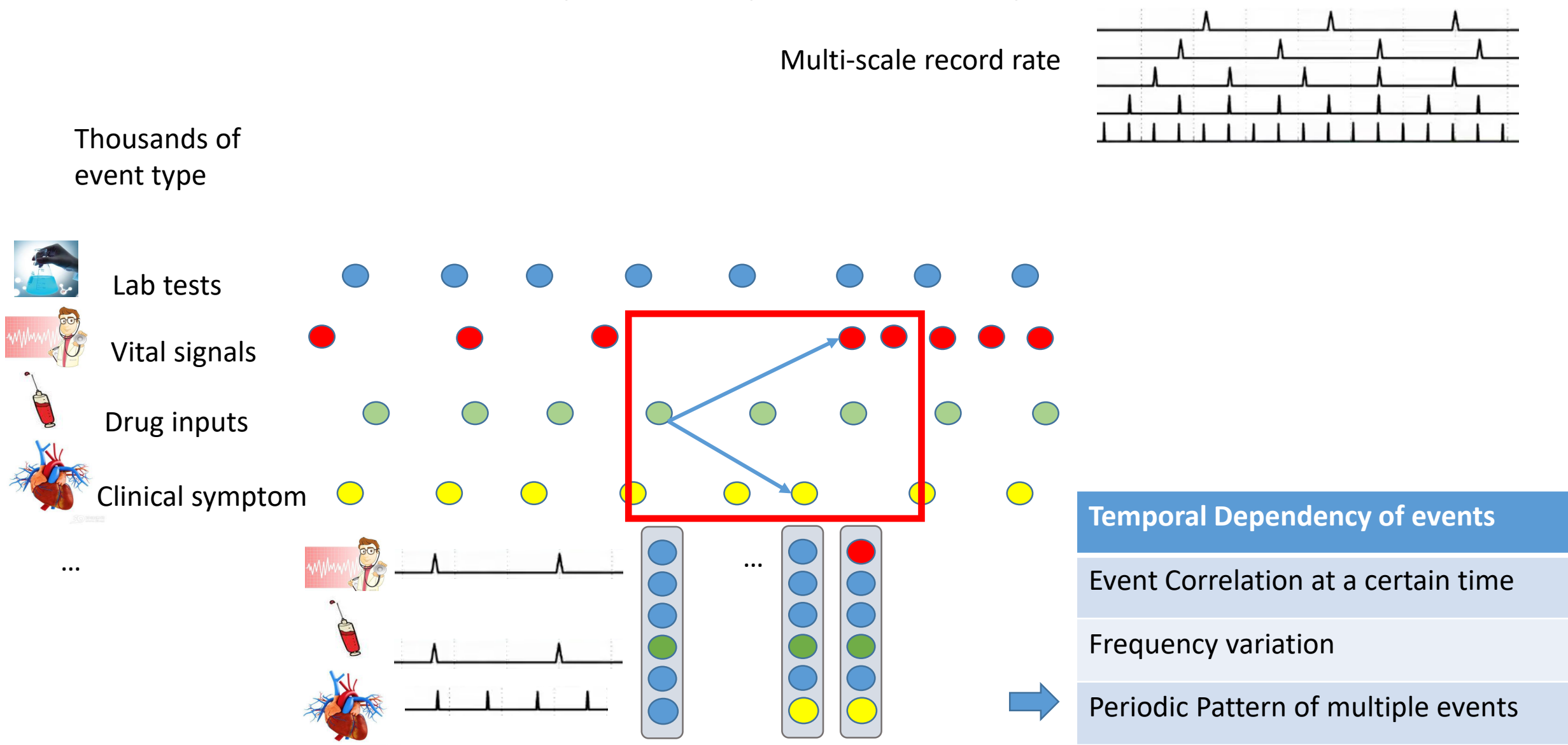


Illustration of the multiple dependency of related events



Outline

- Background—endpoint prediction in EHR
- Heterogeneous temporal events
- Model
- Experiments

Performance on two prediction tasks

Methods	death		lab test	
	AUC	AP	AUC	AP
Independent LSTM	0.8771 ± 0.0005	0.5573 ± 0.0006	0.7196 ± 0.0006	0.2969 ± 0.0008
Independent LSTM(shared weight)	0.8064 ± 0.0005	0.5301 ± 0.0006	0.5308 ± 0.0005	0.1098 ± 0.0005
Phased LSTM	0.8474 ± 0.0005	0.4900 ± 0.0075	0.7722 ± 0.0007	0.3575 ± 0.0026
Clock-work RNN	0.8400 ± 0.0001	0.7181 ± 0.0003	0.6516 ± 0.0002	0.2208 ± 0.0003
RETAIN	0.8967 ± 0.0011	0.5808 ± 0.0114	0.7325 ± 0.0022	0.3096 ± 0.0052
LSTM + event embedding & attr encoding	0.9466 ± 0.0002	0.7445 ± 0.0007	0.7231 ± 0.0028	0.3021 ± 0.0014
HE-LSTM	0.9516 ± 0.0003	0.7687 ± 0.0011	0.7987 ± 0.0008	0.3914 ± 0.0013

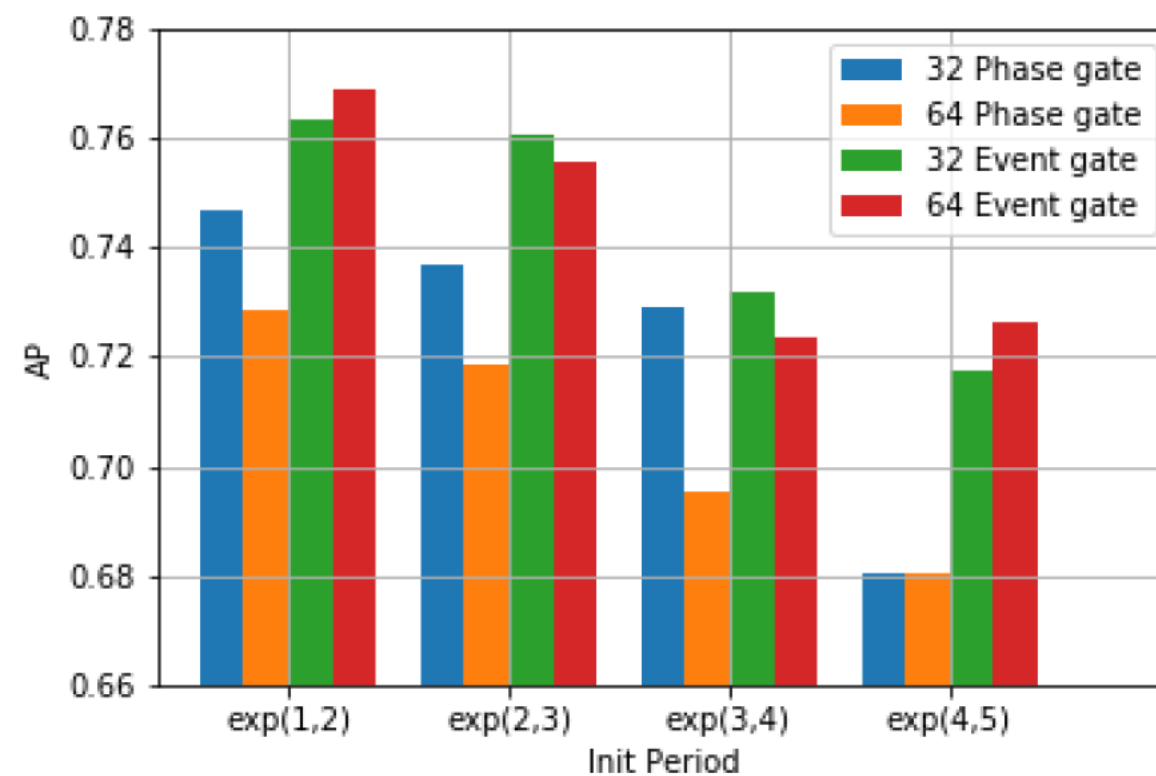
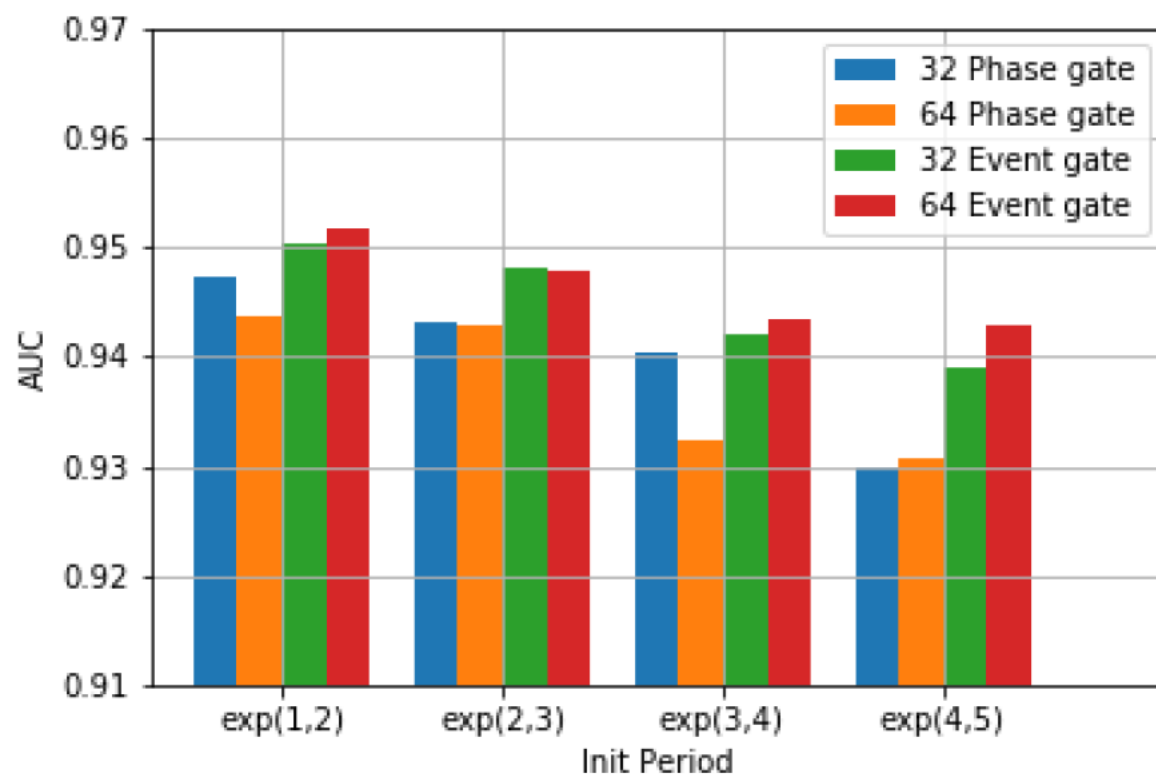
Performance with different settings of event gates

- The phase gate helps to achieve a fast convergence

Methods		Phase gate	Event filter	Event gate
death prediction	AUC(1st epoch)	0.9301	0.9105	0.9370
	AUC	0.9471	0.9518	0.9516
	AP(1st epoch)	0.6856	0.6048	0.7094
	AP	0.7467	0.7679	0.7687
	Entropy(1st epoch)	0.1561	0.1835	0.1479
	Entropy	0.1369	0.1301	0.1297
abnormal lab test prediction	AUC(1st epoch)	0.7050	0.6747	0.7275
	AUC	0.7945	0.7559	0.7987
	AP(1st epoch)	0.2752	0.2403	0.2965
	AP	0.3875	0.3410	0.3914
	Entropy(1st epoch)	0.3373	0.3448	0.3298
	Entropy	0.3019	0.3178	0.3003

Different initial periods

- Robust to different initial period by adding event filter



Conclusion

- The clinical endpoint prediction task based on EHR data
- The representation learning problem of heterogeneous temporal events consists of asynchronous clinical records from multiple sources.
- We propose a novel model called HE-LSTM
 - Modeling the multi-scale sampling rates of different kinds of events and their temporal dependency.



THANK YOU