

# Agenda

- Background
  - healthcare data
  - analytical tasks
  - why deep learning models?
  - deep learning architectures
- **Success of Deep Learning in Computational Healthcare**
  - Medical Classification
  - Sequential Prediction
  - Concept Embedding
  - Data Augmentation
- Open Challenges (10 min)
- Q&A

# Medical Classification

- CAML-NAACL' 18 (free text)
- LEAP (discrete clinical codes)
- Image based Classification (Nature skin cancer)

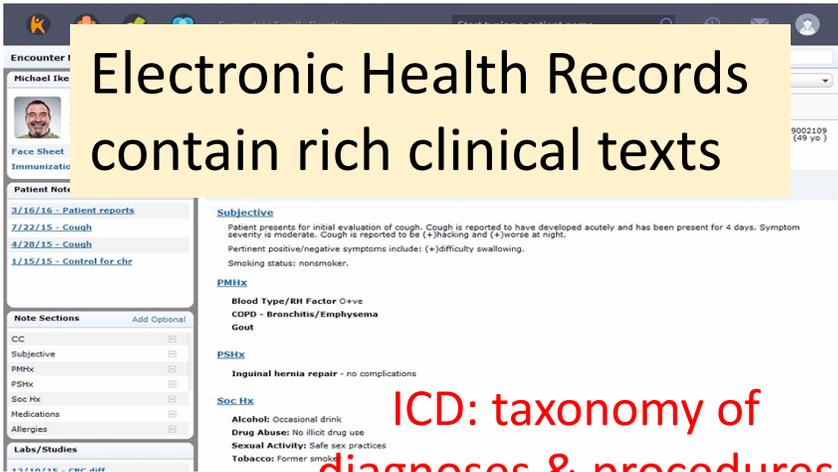
# Explainable Prediction of Medical Codes from Clinical Text

James Mullenbach, Sarah Wiegrefe, Jon Duke, Jimeng Sun, Jacob Eisenstein

**NAACL-HLT' 18**

# CAML: Clinical Coding Problem as Multi-label Classification

## Clinical Coding



Encounter 1  
Michael Ike  
9002109 (49 yo)

Face Sheet  
Immunization

Patient Note

7/16/16 - Patient reports  
7/22/15 - Cough  
9/20/15 - Cough  
1/15/15 - Control for chr

Note Sections Add Optional

CC  
Subjective  
PMHx  
PSHx  
Soc Hx  
Medications  
Allergies

Labs/Studies

Subjective  
Patient presents for initial evaluation of cough. Cough is reported to have developed acutely and has been present for 4 days. Symptom severity is moderate. Cough is reported to be (+)hacking and (+)worse at night.  
Pertinent positive/negative symptoms include: (+)difficulty swallowing.  
Smoking status: nonsmoker.

PMHx  
Blood Type/RH Factor O+ve  
COPD - Bronchitis/Emphysema  
Gout

PSHx  
Inguinal hernia repair - no complications

Soc Hx  
Alcohol: Occasional drink  
Drug Abuse: No illicit drug use  
Sexual Activity: Safe sex practices  
Tobacco: Former smoker

Electronic Health Records contain rich clinical texts

ICD: taxonomy of diagnoses & procedures



## Multi-label Classification

- Highly multi-label classification
  - 14K ICD9, 68K ICD10 labels
- Testbed for document representations
- Documents are long and loosely structured

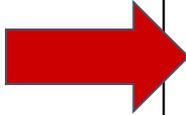
Features	ICD-9
Possible Codes	14,000
Characters	3-5

Laborious and error-prone

# CAML: The MIMIC-III Dataset

- Open-access, de-identified
- 47k admissions -> 47k documents for training
- **Loosely structured:**

## Many labels:

Admission Date: <code>[**2118-6-2**]</code>	Discharge Date: <code>[**2118-6-14**]</code>	519.1: 'Other disease..'
Date of Birth:	Sex: F	491.21: 'Obstructive ...'
Service: MICU and then to <code>[**Doctor Last Name **]</code> Medicine		518.81: 'Acute respir...'
HISTORY OF PRESENT ILLNESS: This is an 81-year-old female with a history of emphysema (not on home O2), who presents...		486: 'Pneumonia, orga...'
		276.1: 'Hyposmolality...'
		244.9: 'Unspecified h...'
		31.99: 'Other operati...'
		.
		.
		.

**Long:** Median post-processed document length: 1,341

Median # labels: 14

# CAML: Modeling Consideration

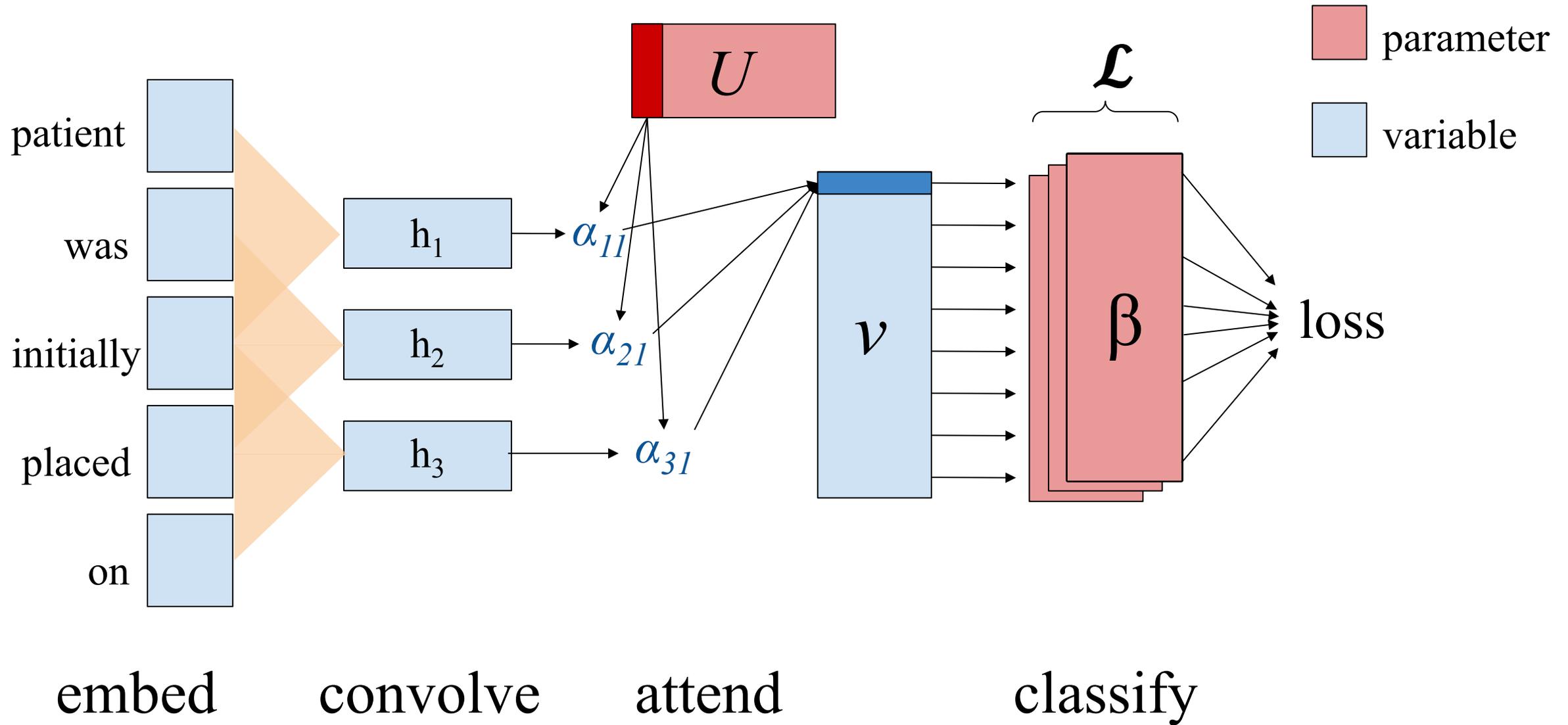
- Focus on the parts that matter
- Treat labels *individually*
- Be fast!

E849.0: Home accidents

801.26: ...subdural,  
and extradural  
hemorrhage...

...who sustained **a fall at home** she was found to  
have a large acute on **chronic subdural hematoma**  
with extensive midline shift...

# The CAML model



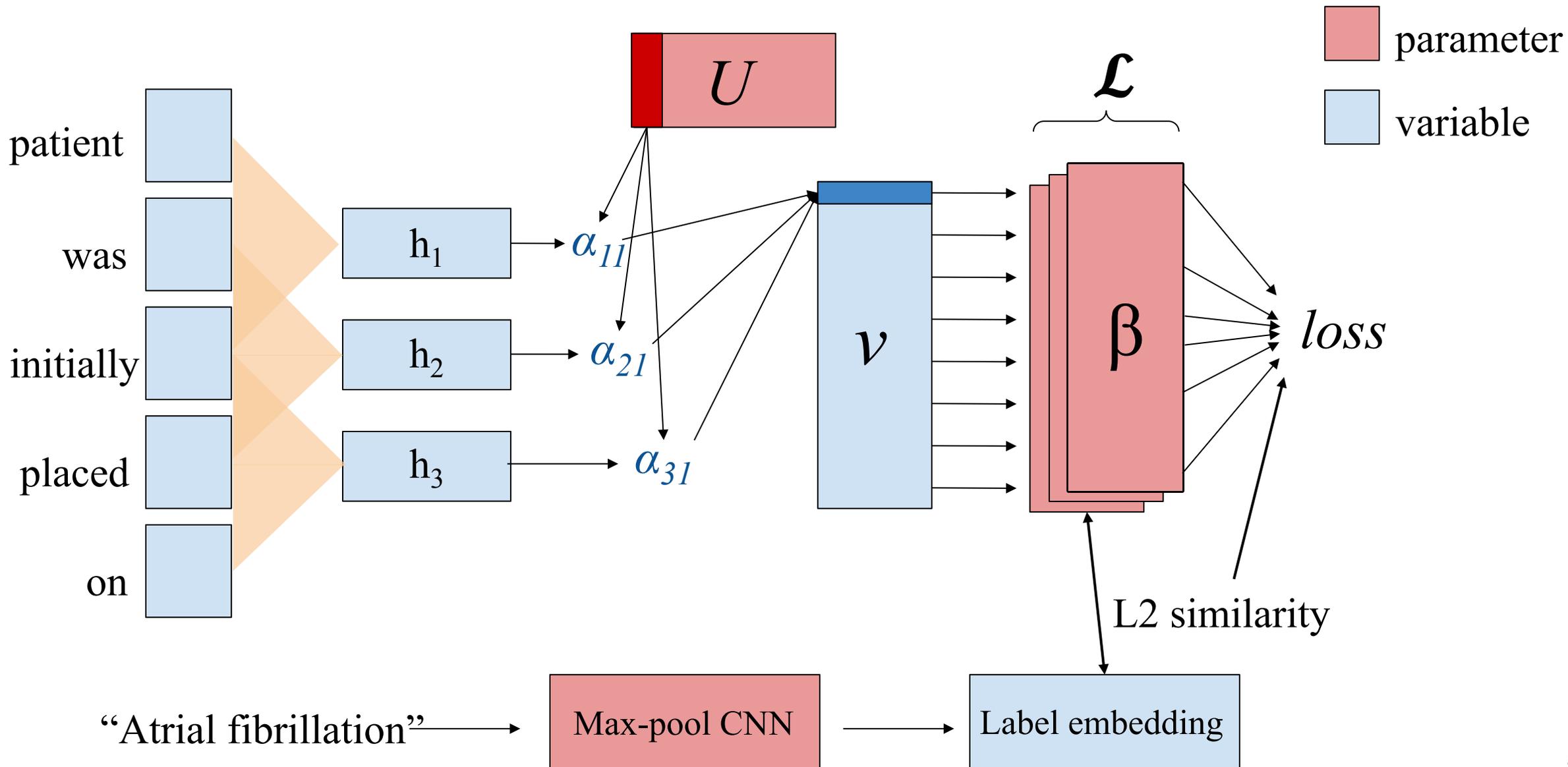
# CAML: Dealing with the Long Tail

- Huge label space (nearly 9,000 total)
- Many labels are similar

250.00: “Diabetes mellitus without mention of complication, type II or unspecified type, not stated as uncontrolled”

250.02: “Diabetes mellitus without mention of complication, type II or unspecified type, uncontrolled”

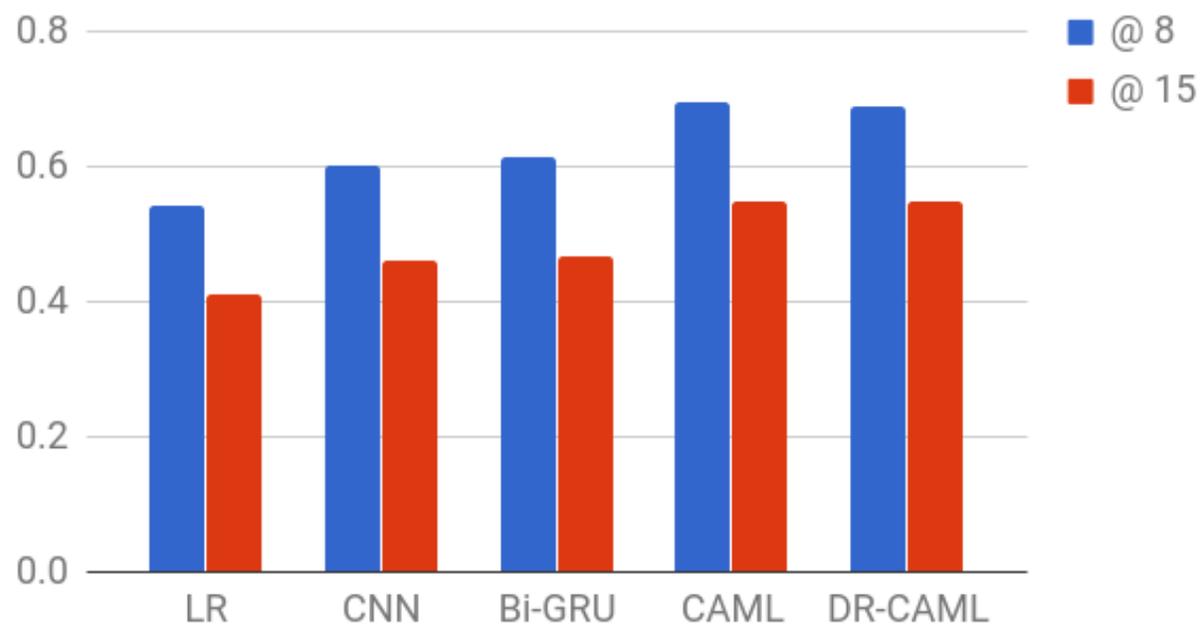
# DR-CAML



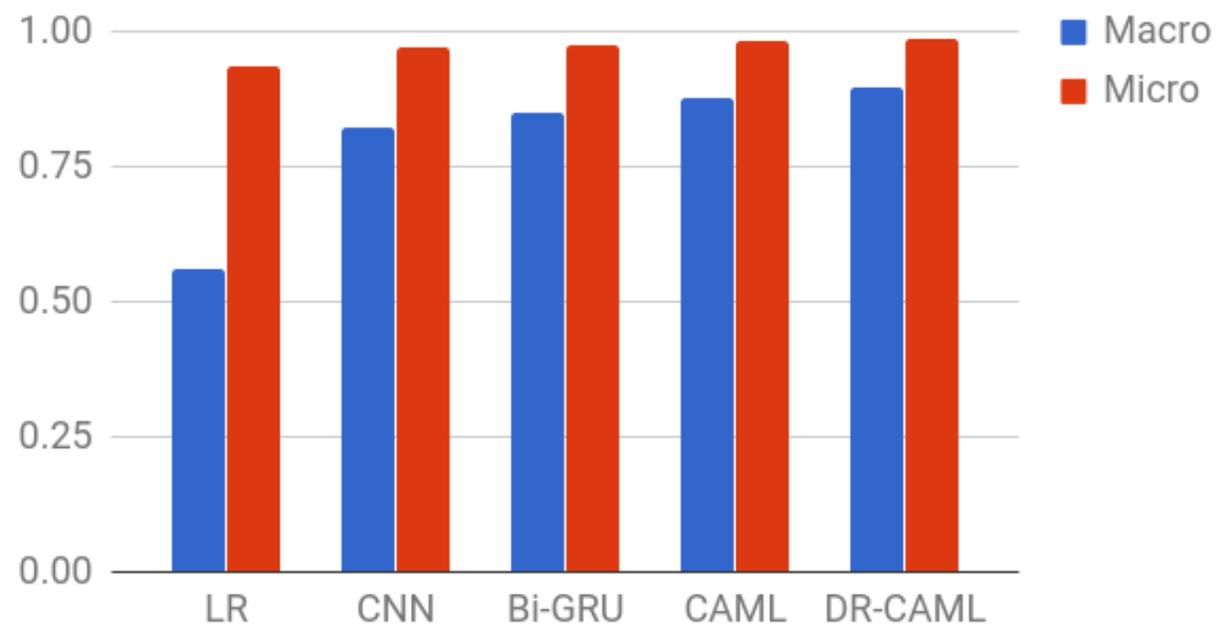
# CAML: Experiment Results

- Enable future comparison
- Precision @ k: decision support use-case

Precision @ k

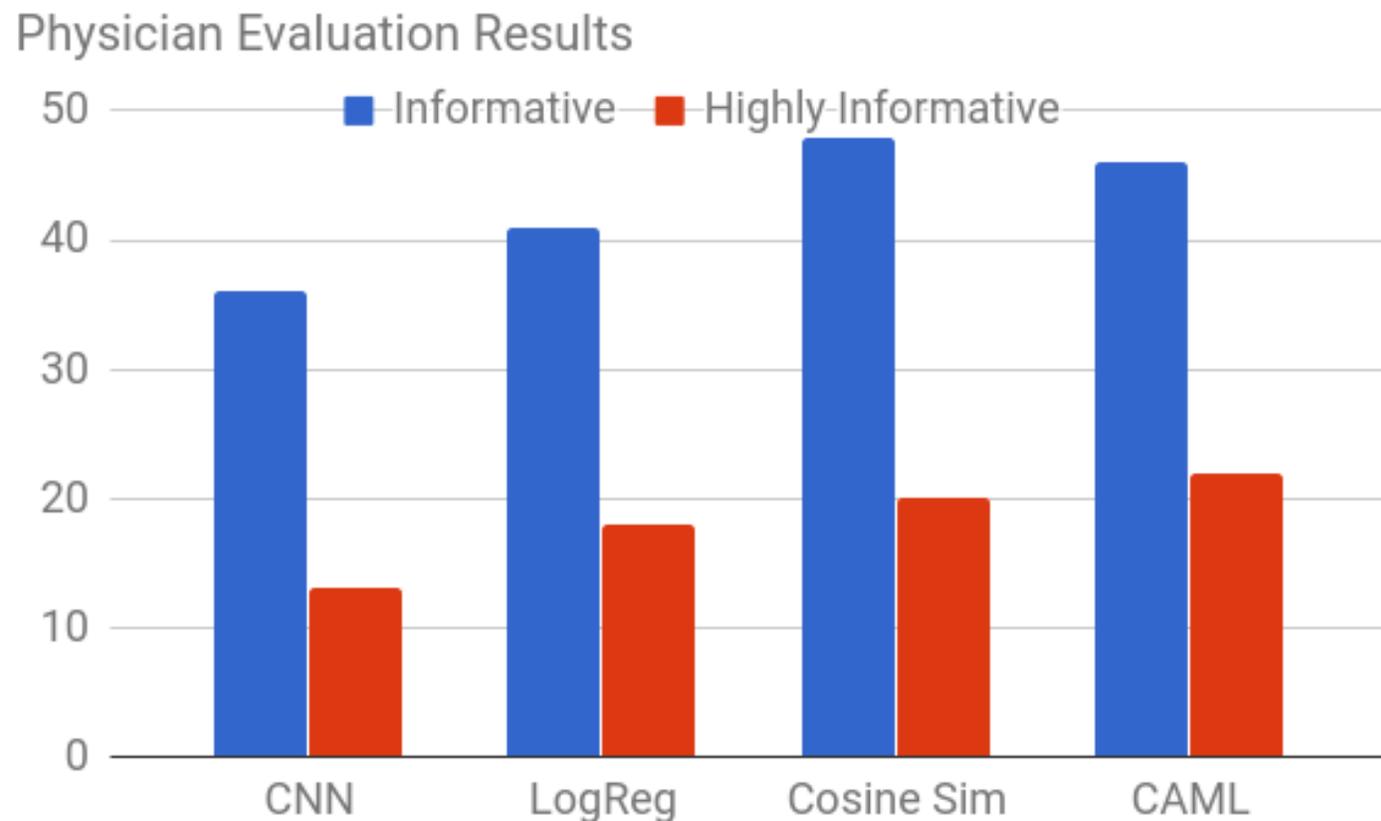


Macro- and Micro-AUC



# CAML: Physician Evaluation Results

- Improves upon CNN, LogReg
- More experts needed!



# CAML: Physician Evaluation Example

Code: 575.4

Full descriptions: Perforation of gallbladder

“. . . in the setting of gallbladder perforation secondary to acute acalculous cholecystitis after . . . . . inhalation hospital 1 times a day metronidazole mg tablet sig one tablet po tid times . . . . . to have an infection in your gallbladder requiring iv antibiotics and tube placement for . . . . .”

LogReg

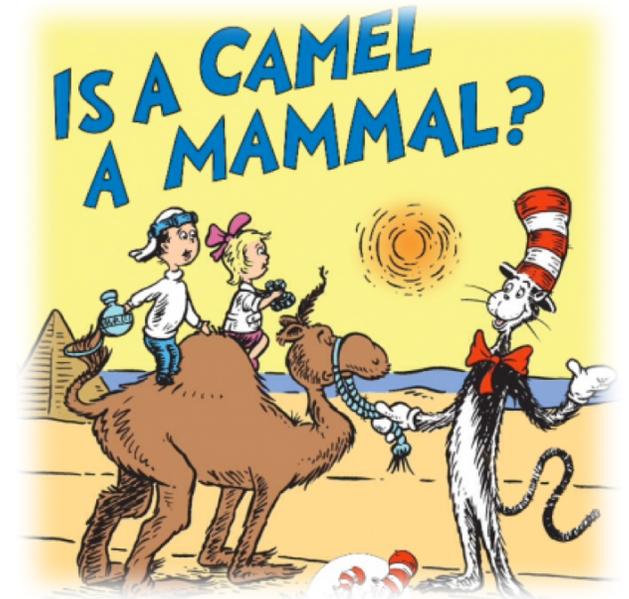
CAML

cosine sim.

CNN

# CAML: Summary

- ICD coding is valuable and challenging
- Convolution + attention works well
- Attention can explain the predictions



# LEAP: Learning to Prescribe Effective and Safe Treatment Combinations for Multimorbidity

Yutao Zhang, Robert Chen, Jie Tang, Walter F. Stewart, Jimeng Sun

**KDD' 17**

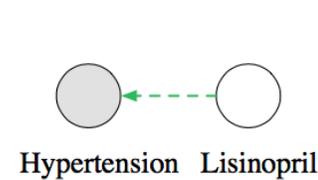
# LEAP: Multi-morbidity

- Co-occurrence of multiple medical conditions
- Traditional way of prescribing is based on doctors' intuition.
- Clinical decisions can be sub-optimal due to knowledge gaps.

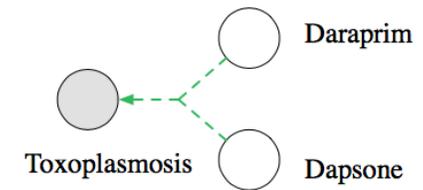


# Complex Dependency between Drugs and Diseases

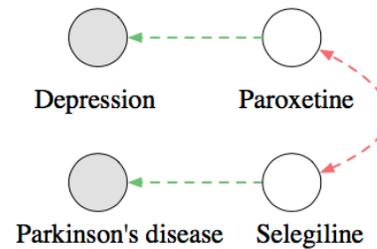
- One drug can treat multiple diseases.
- One disease requires a combination of multiple drugs.
- Potential adverse interactions among drugs.
- Potential conflicts between drugs and diseases.



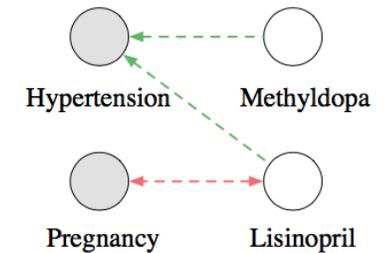
(a) One to one mapping



(b) Many to one mapping



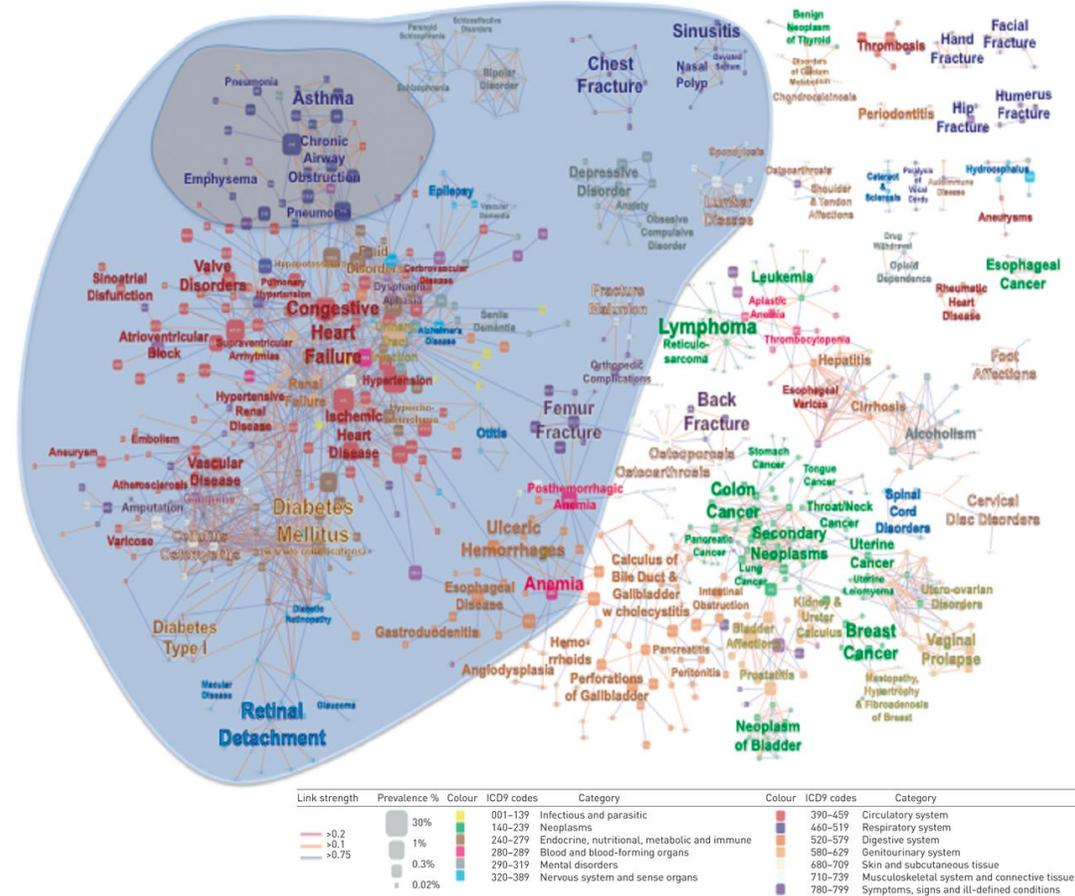
(c) Drug-drug interaction



(d) Drug-disease interaction

# LEAP: Treatment Recommendation for Multi-morbidity Patients

- Many-to-many mapping:
  - Patients are diagnosed with multiple diseases.
  - Requires a combination of drugs and a treatment.
- Learning from EHR data:
  - No explicit mapping between drugs and diseases.
  - Complex dependencies.



# LEAP: Multi-Instance Multi-Label Learning

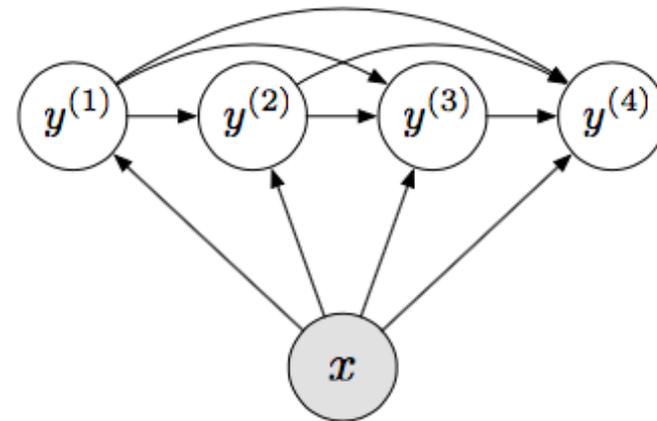
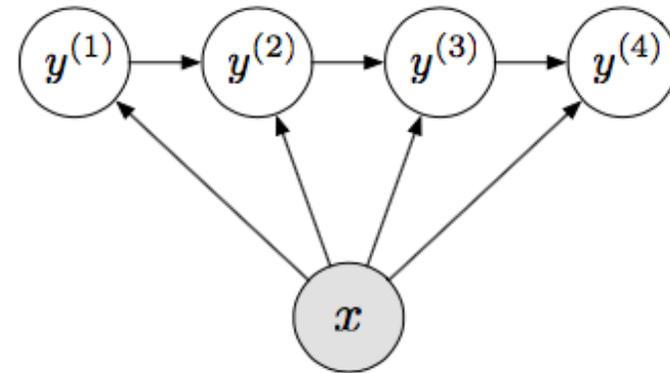
- Multi-Instance Learning:
  - A discrete set is given as input.
  - Assign a label to each input set.
- Multi-Label Learning:
  - Predict a variable length label set for each instance.

# LEAP: Multi-Instance Learning Problem

- Independence assumption
  - The instances within an input set are independent.
- Equal contribution assumption
  - Each instances within an input set contribute equally to the generation of the corresponding label.
- Not true in the setting of treatment recommendation due to the complex high-order dependencies.

# LEAP: Multi-Label Learning Problem

- Binary Relevance
  - The simplest solution.
  - Assume independency among labels.
  - Translate the Multi-label problem into a set of binary problems.
- Classifier Chains
  - Use the previous predictions as input of the next classifier.



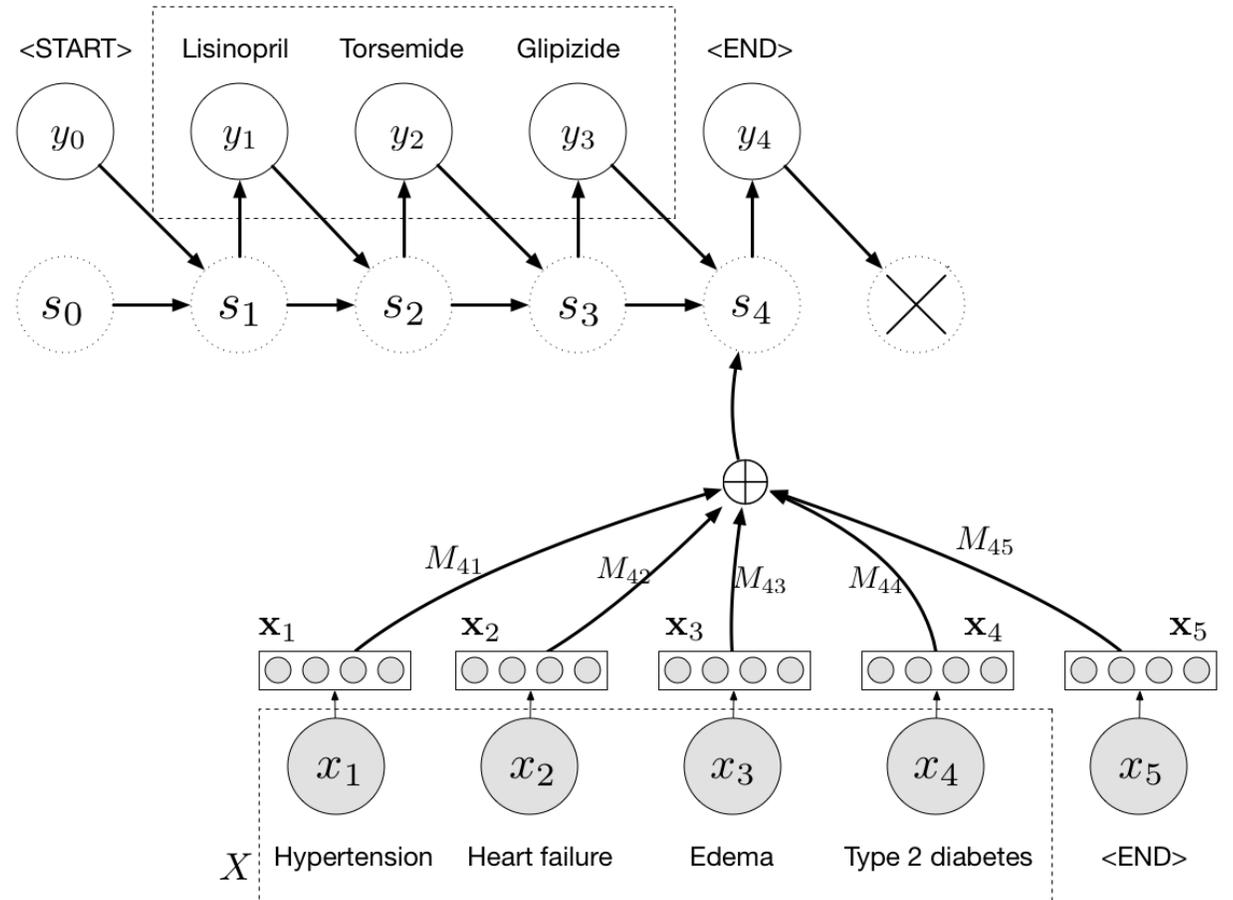
# LEAP Model

- Objective: Conditional probability  $p(Y|X)$  of drug set  $X$  given disease set  $Y$
- Inspired by Classifier Chains:
  - Decompose the **Combinatorial Optimization** problem as a **Sequential Decision Making** problem:

$$\begin{aligned} p(Y|X) &= \prod_{t=1}^{|Y|} p(y_t | X, y_1, y_2, \dots, y_{t-1}) \\ &= \prod_{t=1}^{|Y|} p(y_t | \{x_1, x_2, \dots, x_{|X|}\}, y_1, y_2, \dots, y_{t-1}), \end{aligned}$$

# LEAP Model

- At step  $t$ , predict the next label based on input set  $X$  and the previous predictions.
- Different instances in the input set contribute differently to the currently prediction.
  - Eg. The weight of an input disease should be decreased if it has already been covered by some predicted drug.

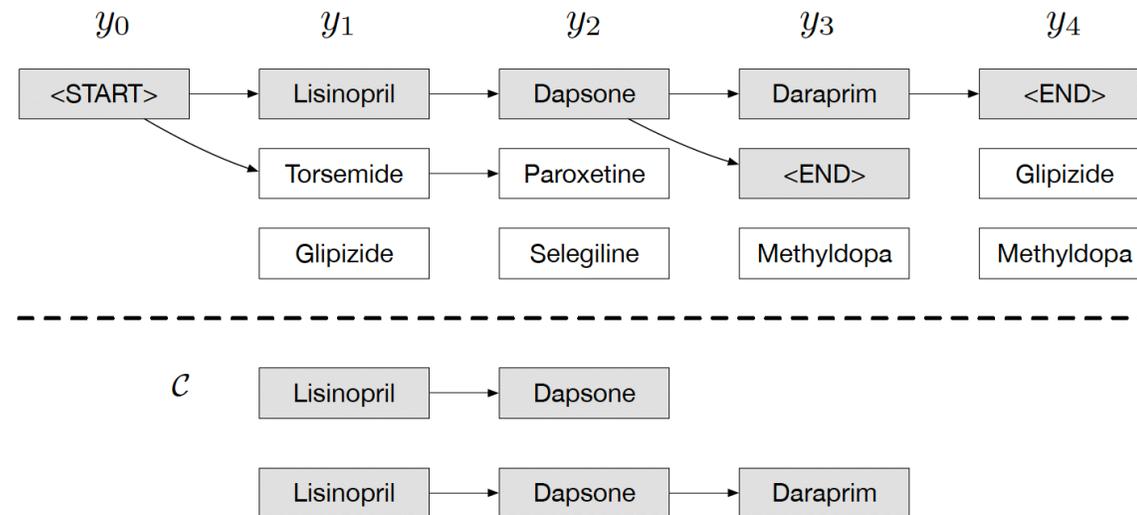


# LEAP Model

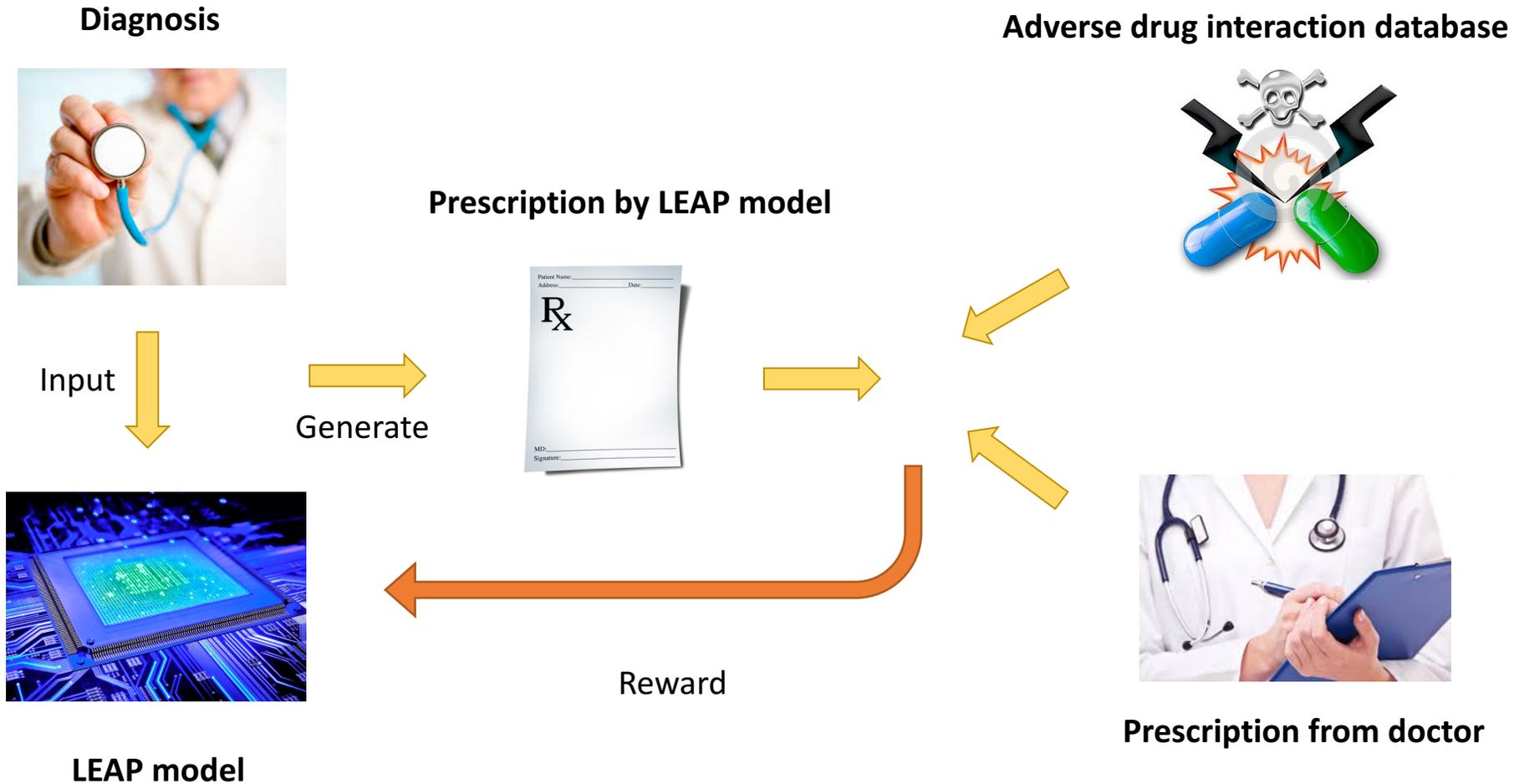
- $s_t = g(s_{t-1}, y_{t-1}, \Psi_t(X))$  the state variable at step t.
  - state : variable at step t-1  $s_{t-1}$
  - $y_{t-1}$  prediction at step t-1
  - $\Psi_t(X)$  an attention function over the input set, capturing the compatibility between input instances and state variable
- Use RNN to model the Sequential Decision Making process.

# LEAP: Beam Search Decoding

- Starting from a <START> label.
- Keep the top-K prediction results.
- Move the prediction path into the candidate set when a <END> label is predicted.
- Terminates if there is no better prediction path.



# LEAP: Reinforcement Fine-Tuning



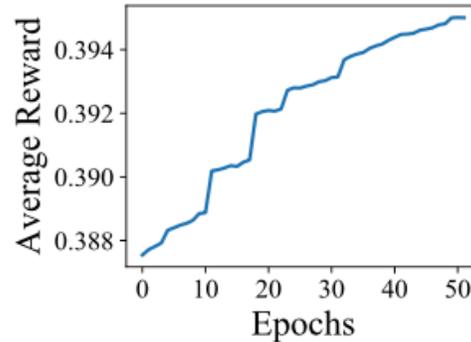
# LEAP: Experimental Results

- Evaluated on real world EHR data from Sutter(heart disease and epilepsy patients) and MIMIC-3(ICU patients).
- Use the first and third level of GPI codes and the candidate drug set.

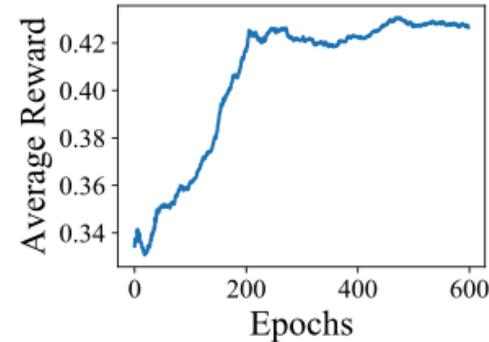
	<b>Sutter</b>		<b>MIMIC-3</b>	
<b>Granularity</b>	<b>1</b>	<b>3</b>	<b>1</b>	<b>3</b>
Rule-based	0.3207	0.2770	0.2753	0.2354
<i>K</i> -Most frequent	0.4283	0.3181	0.2609	0.2616
Softmax MLP	0.4908	0.3739	0.4897	0.3342
Classifier Chains	0.4839	0.3620	0.4621	0.3204
Basic LEAP	0.5270	0.3936	0.5107	0.3865
LEAP	<b>0.5341</b>	<b>0.4073</b>	<b>0.5582</b>	<b>0.4342</b>

# LEAP: Reinforcement Fine-tuning

- Average reward w.r.t training epochs



(a) Sutter Dataset



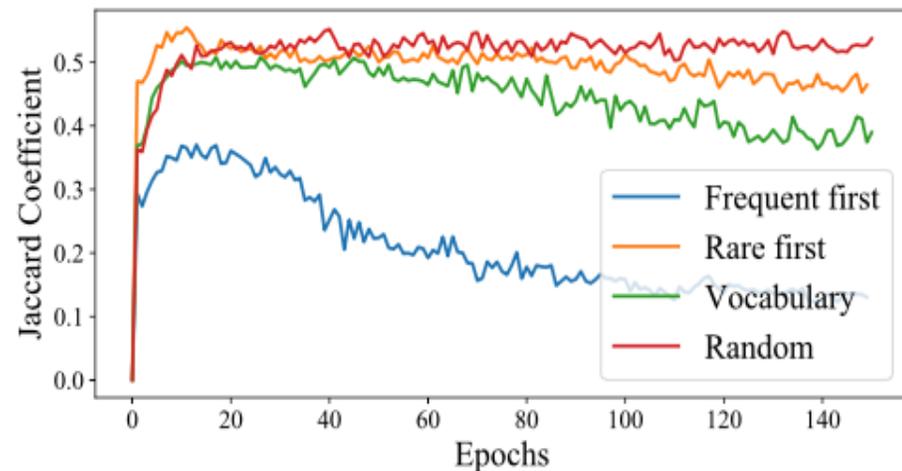
(b) MIMIC-3 Dataset

- Performance of avoiding adverse drug interactions.

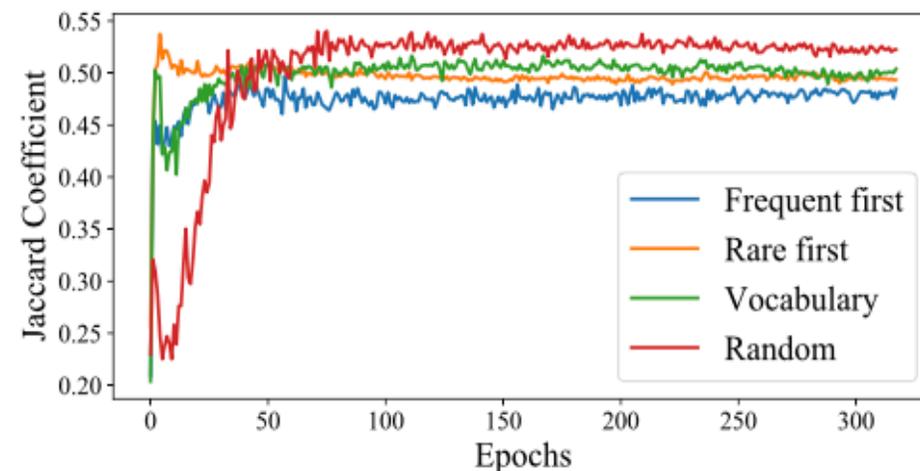
Method	Drug Interaction Rate
<i>K</i> -Most frequent	12.06%
Softmax MLP	3.51%
Basic LEAP	2.41%
LEAP	0.23%

# LEAP: Different Order of Labels

- Frequent-first performs the worst.
- Rare-first converges fast.
- Random-shuffle is the most robust one.



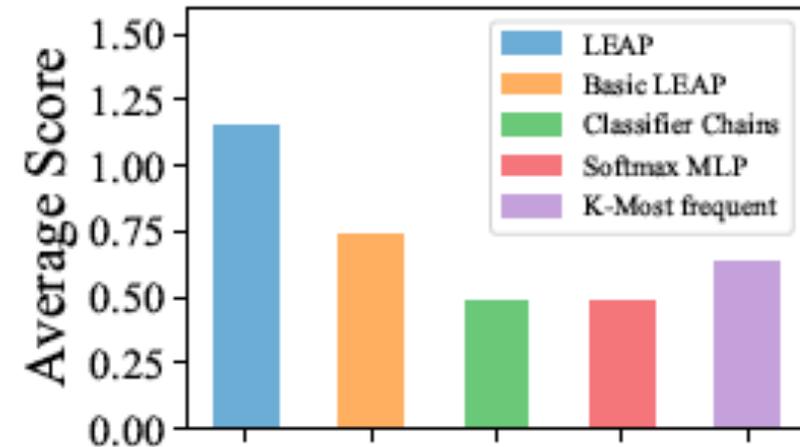
(a) Sutter Dataset



(b) MIMIC-3 Dataset

# LEAP: Qualitative Evaluation

- The results are manually scored by a clinical expert
  - 2: Completely covers the diseases with no conflict.
  - 1: Partially covers the diseases (>50%) with no conflict.
  - 0: Conflict or incomplete (<50%)
- Case study:



Diagnosis	Methods	Recommended Treatments
Type 2 diabetes	K-Most frequent	PEG KCl Bicarb, Quinapril, Pravastatin, Metformin, Paroxetine
Hyperlipidemia	Softmax MLP	Metformin
Depressive disorder	Basic LEAP	Metformin, Quinapril, Pravastatin, Paroxetine
Hypertension	LEAP	Metformin, Amiloride/HCTZ, Fenofibrate, Paroxetine
Depressive Disorder	K-Most frequent	Azithromycine, Privastatin, Paroxetine
Acute bronchitis	Softmax MLP	Paroxetine, Azithromycin
Imbalance (gait)	Basic LEAP	Azithromycin
	LEAP	Azithromycin, Dextromethorphan-K Guaiacolsulfonate, Paroxetine

**Table 3: Example Recommended Treatments by Different Methods**

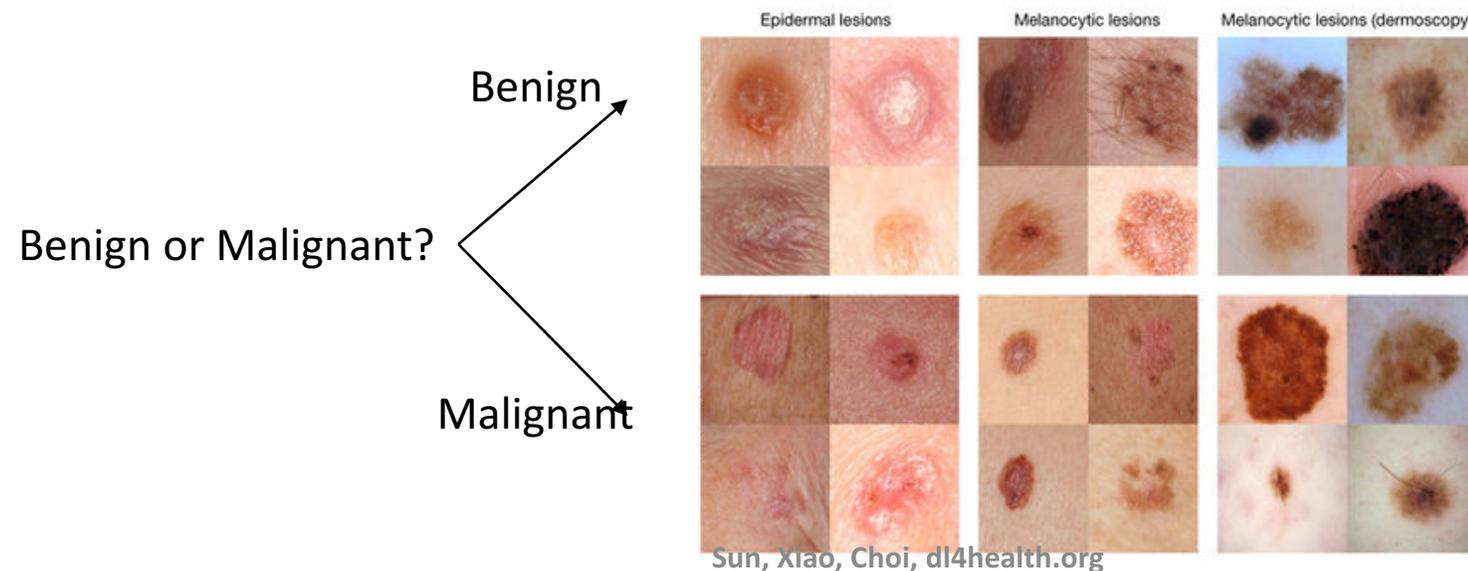
# Dermatologist-level classification of skin cancer with deep neural networks

A. Esteva, B. Kuprel, R.A. Novoa, J. Ko, S.M. Swetter, H.M. Blau, S. Thrun

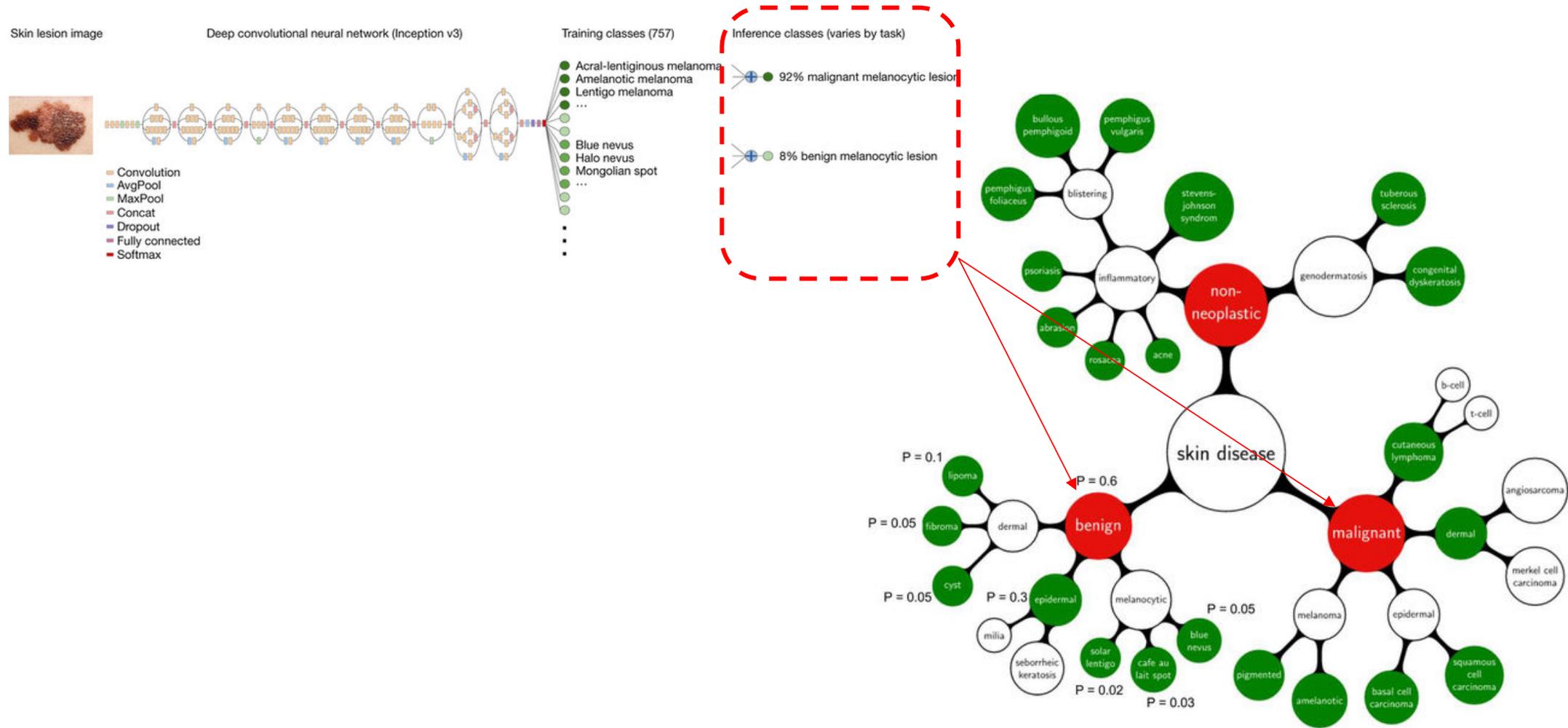
**Nature 542, p115-118 (2017)**

# Detecting Skin Cancer

- Dermatologist-level classification of skin cancer with deep neural networks
  - Esteva, Andre, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun, Nature 2017
- Given clinical images, classify
  - Keratinocyte carcinomas VS benign seborrheic keratoses
  - malignant melanomas VS benign

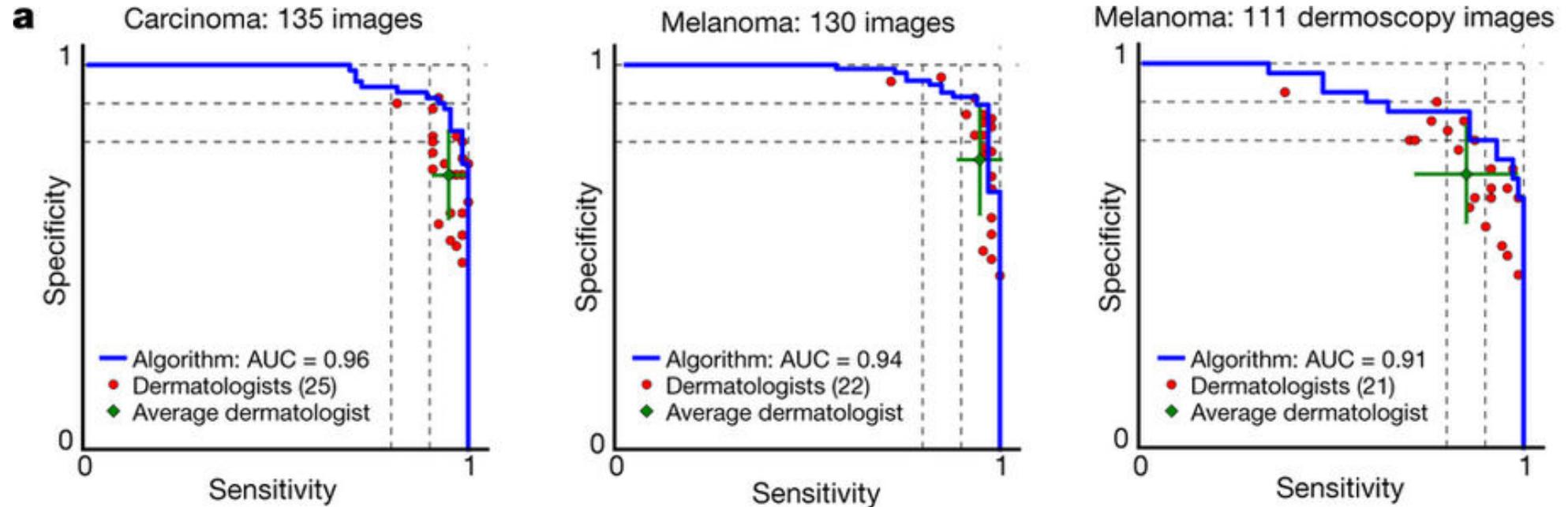


# Model architecture: inception v3



# Model performance

- Better than average board-certified dermatologists



# Sequential Prediction

- Dr. AI/JAMIA-HF-RNN (EHR- multilabel-binary)
- RETAIN (EHR-interpretability)
- RAIM (multimodal)
- CONTENT (readmission prediction, PLOS One)

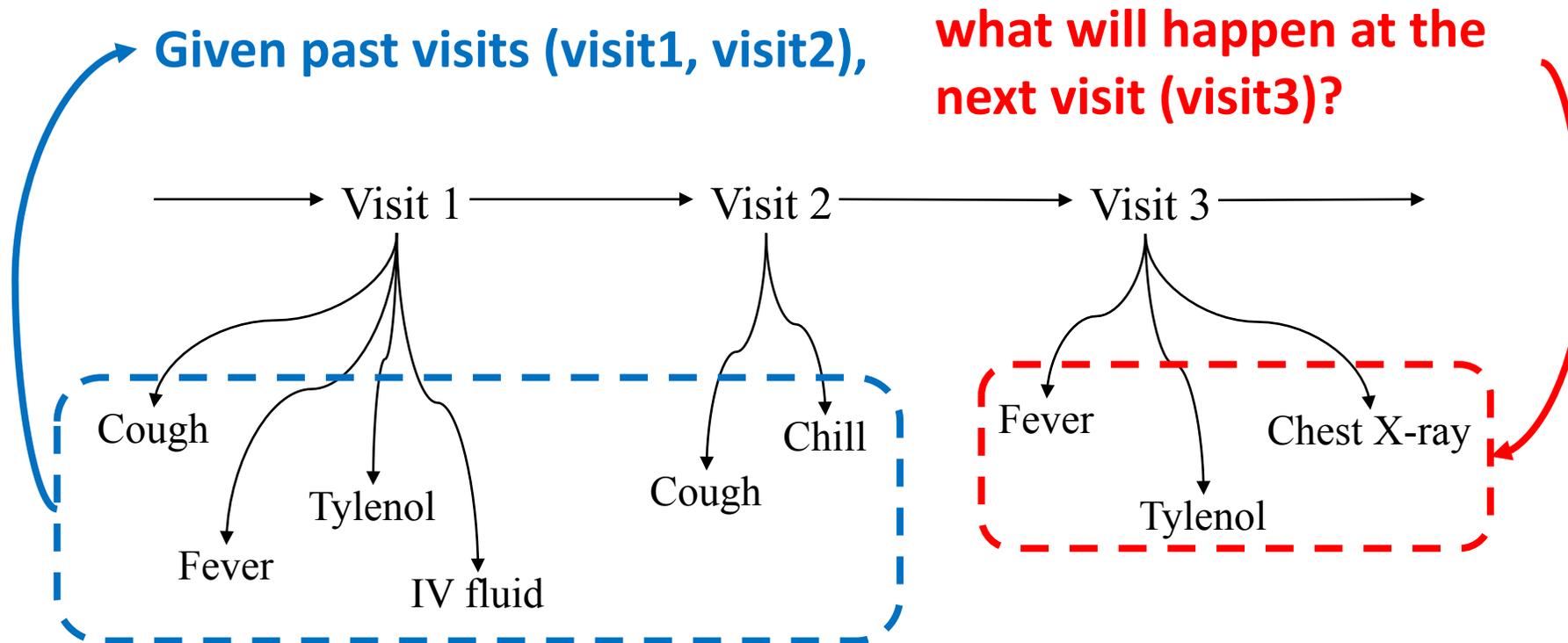
# Doctor AI: Predicting Clinical Events via Recurrent Neural Networks

Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F. Stewart, Jimeng Sun

*Machine Learning for Healthcare Conference, 2016*

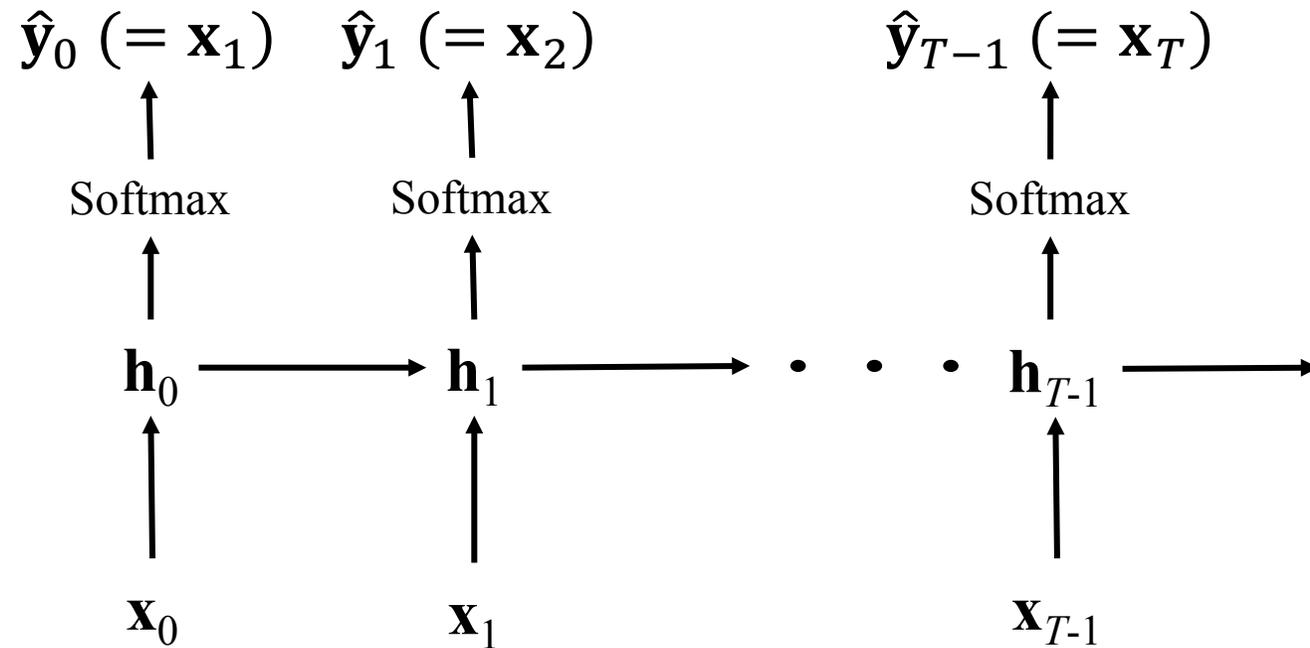
# Doctor AI: Background

- Disease progression modeling



# Doctor AI: Model

- Feed visits into the RNN
  - One visit at each timestep.
  - Predict next events at each timestep.

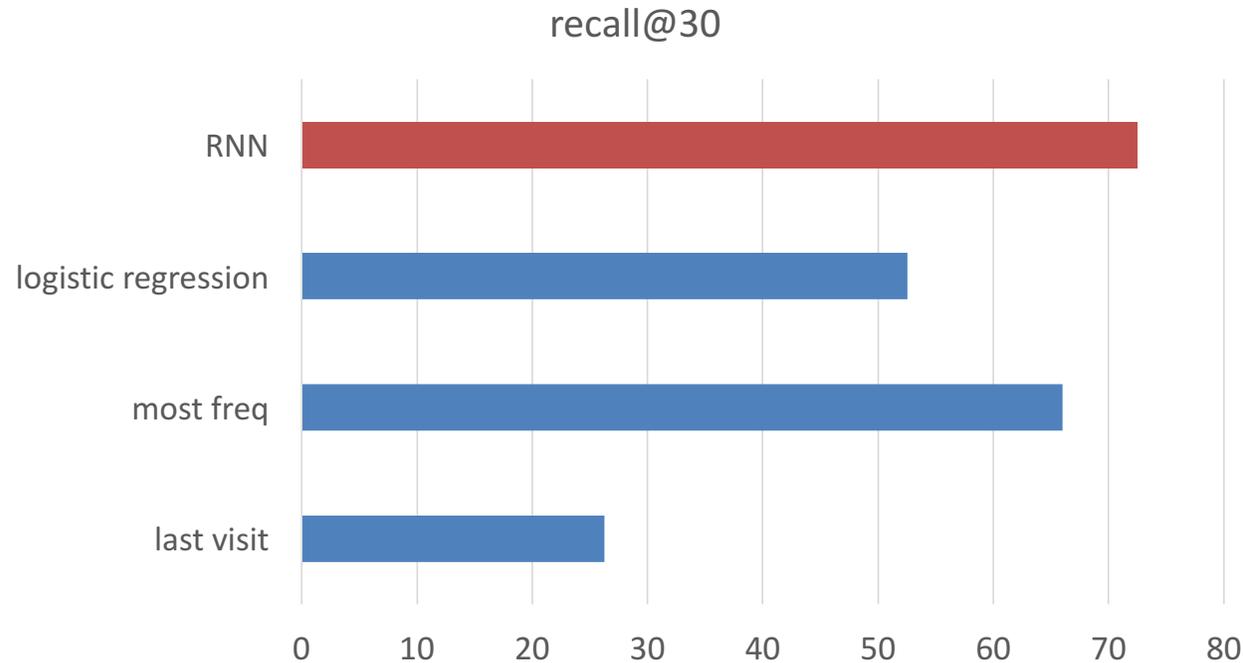


# Doctor AI: Data

- Data
  - 260K patients from Sutter Health
  - Patient records over 10 years
  - Input codes
    - Diagnosis codes, medication codes, procedure codes (38,000 codes)
  - Output labels
    - 1,183 diagnosis codes

# Doctor AI: Sequential Prediction

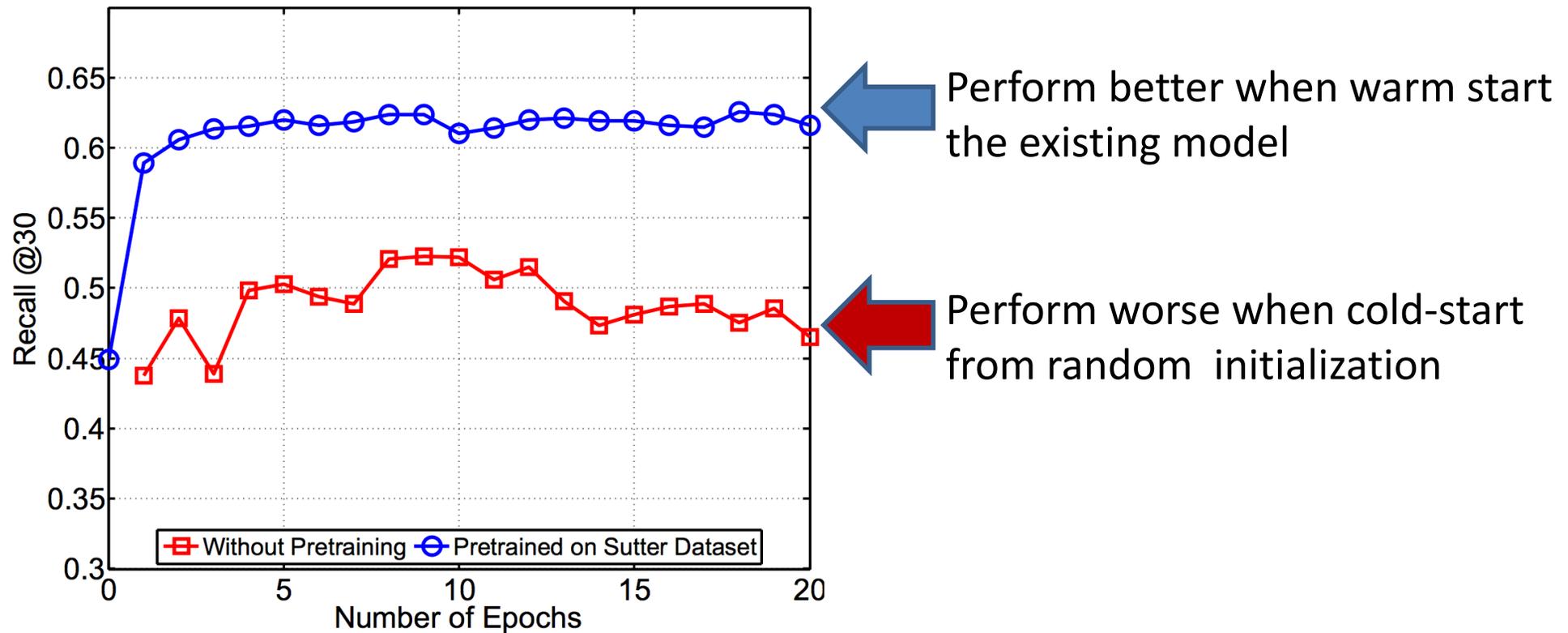
- Predicting diagnoses in the next visit



$$\text{top-}k \text{ recall} = \frac{\# \text{ of true positives in the top } k \text{ predictions}}{\# \text{ of true positives}}$$

# Doctor AI: Knowledge Transfer

- Generalize RNN model from one hospital to another



# RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism

Edward Choi, Mohammad Taha Bahadori, Joshua A. Kulas, Andy Schuetz, Walter F. Stewart, Jimeng Sun

**NIPS' 16**

# RETAIN: Interpretable Models

What do we mean by interpretability of a model?

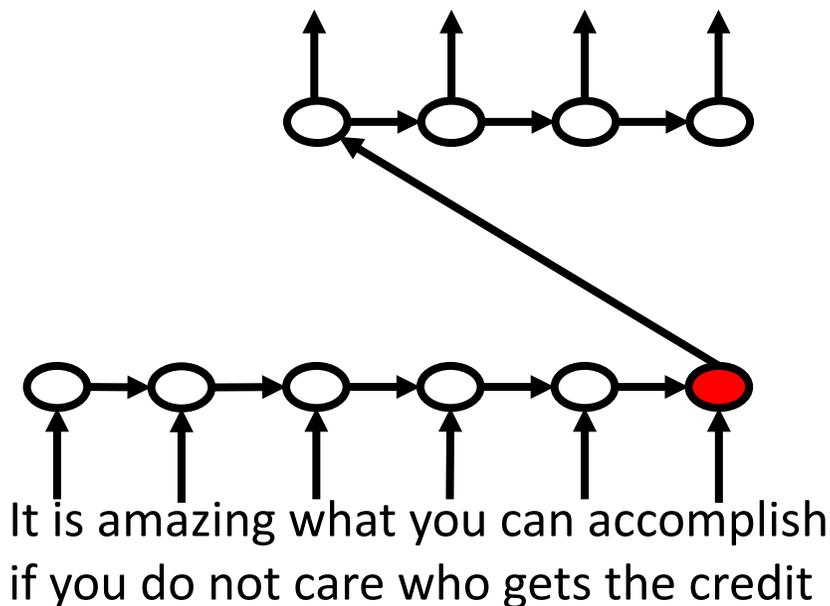
## Three categories of models:

- Rule based: *e.g.* decision trees
  - Case based: *e.g.* nearest neighbor methods
  - Risk factor based: *e.g.* sparse linear regression
- 
- Temporal models? Latent variable models

# RETAIN: Neural Attention Mechanism

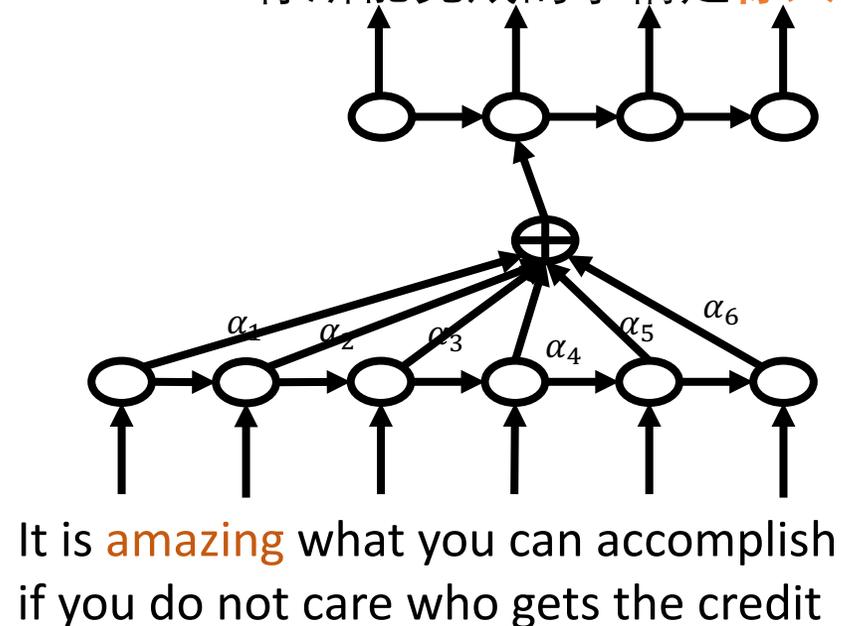
## Regular Machine Translation

如果你不在乎谁获得了荣誉，  
你能完成的事情是惊人的。



## Neural Attention Mechanism

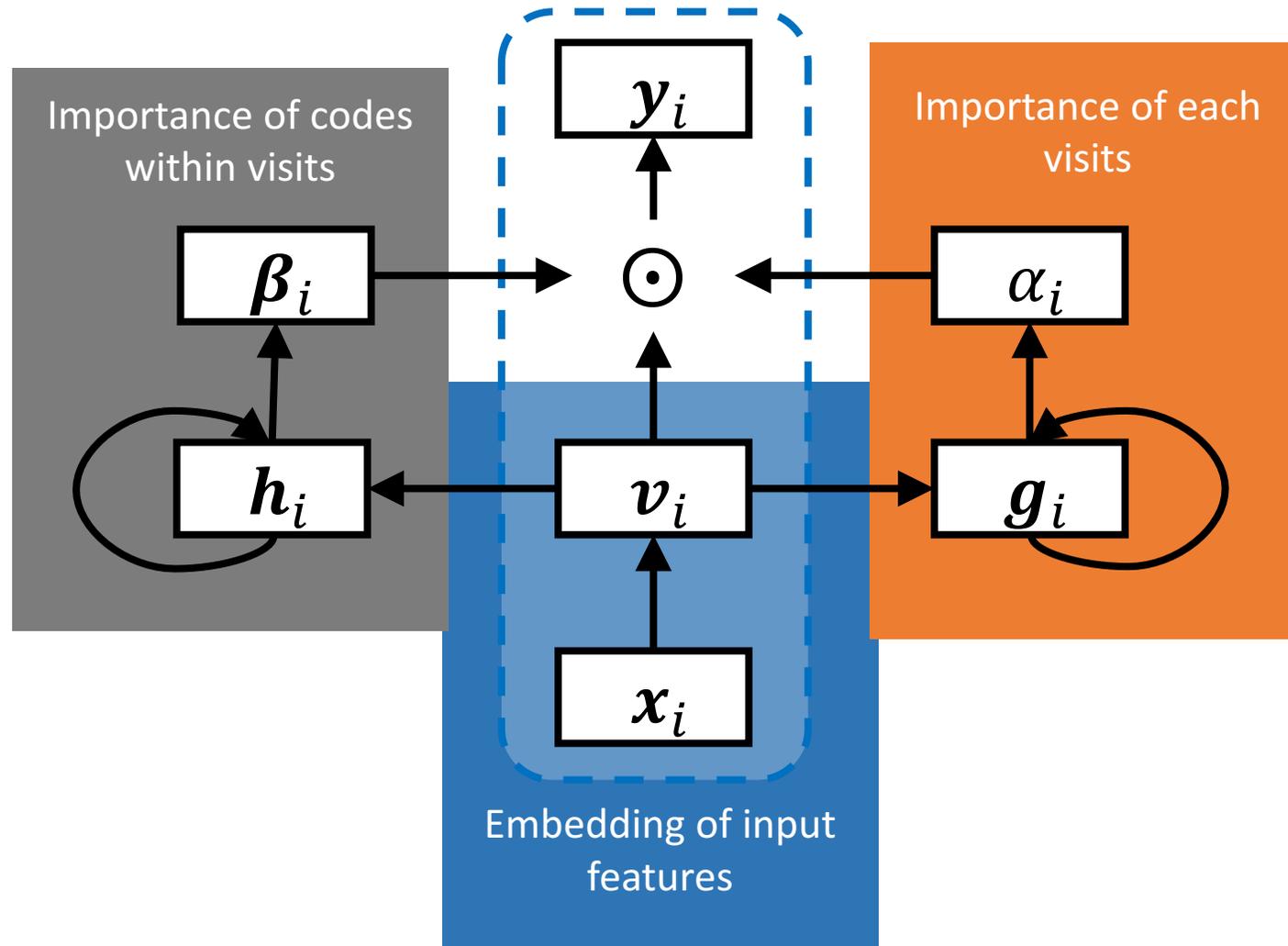
如果你不在乎谁获得了荣誉，  
你能完成的事情是**惊人的**。



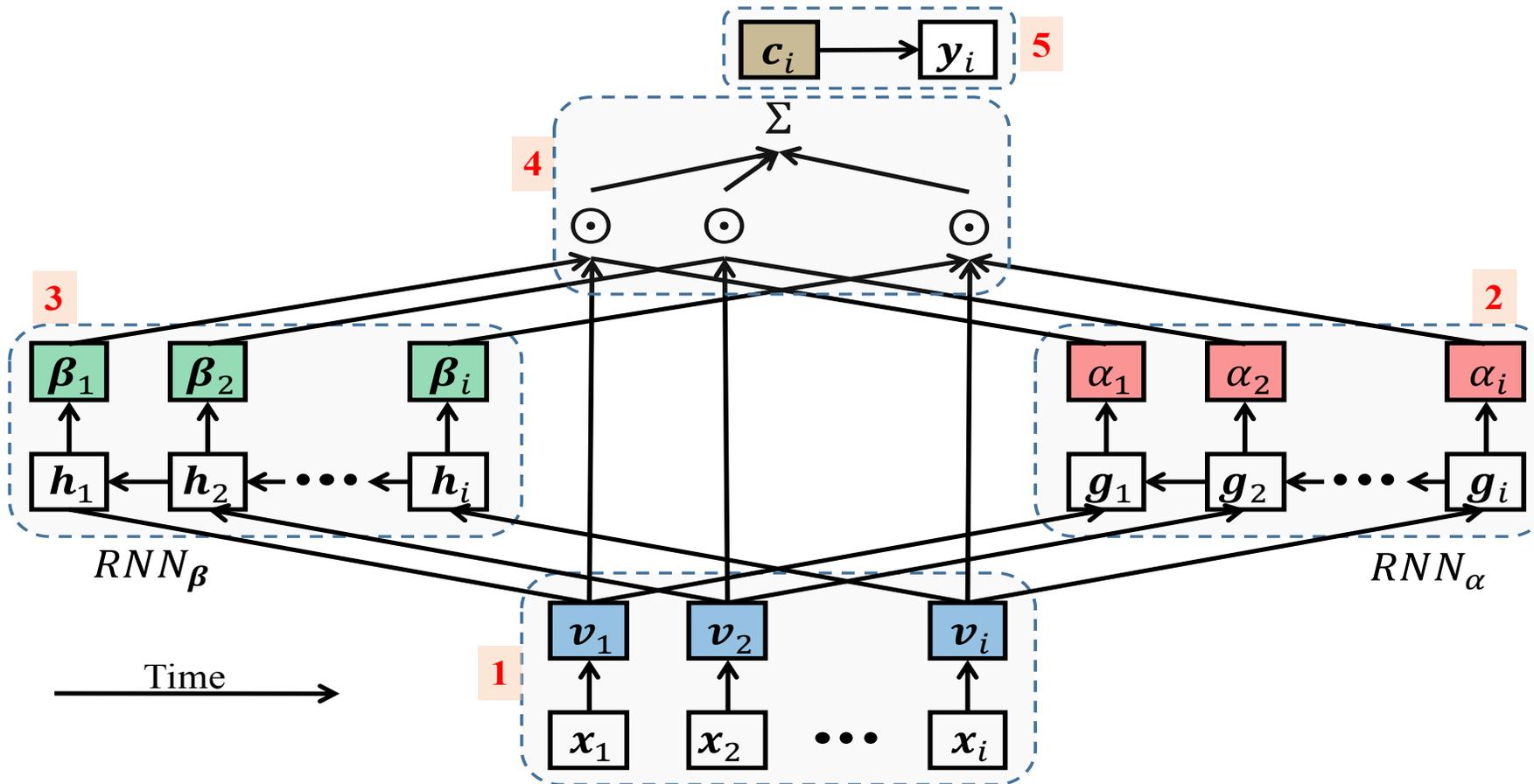
Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2014. "Neural Machine Translation by Jointly Learning to Align and Translate."

Sun, Xiao, Choi, dl4health.org

# RETAIN: REverse Time Attention model

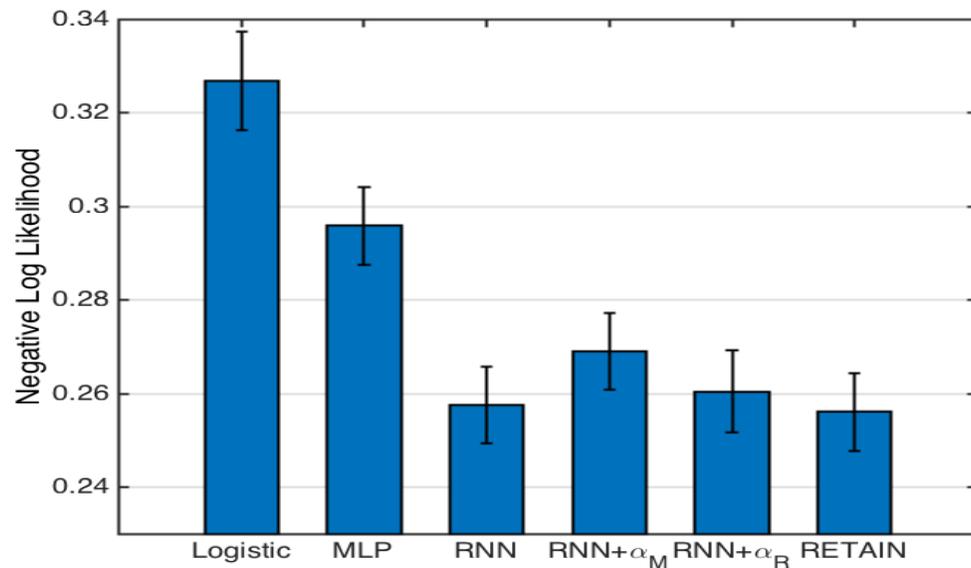


# Details of RETAIN

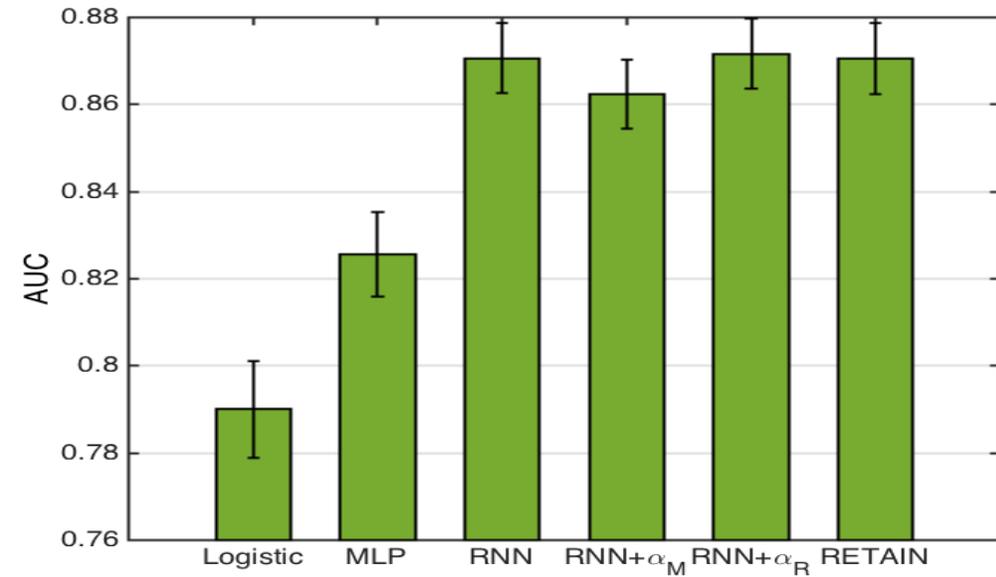


# RETAIN: Heart Failure Prediction Results

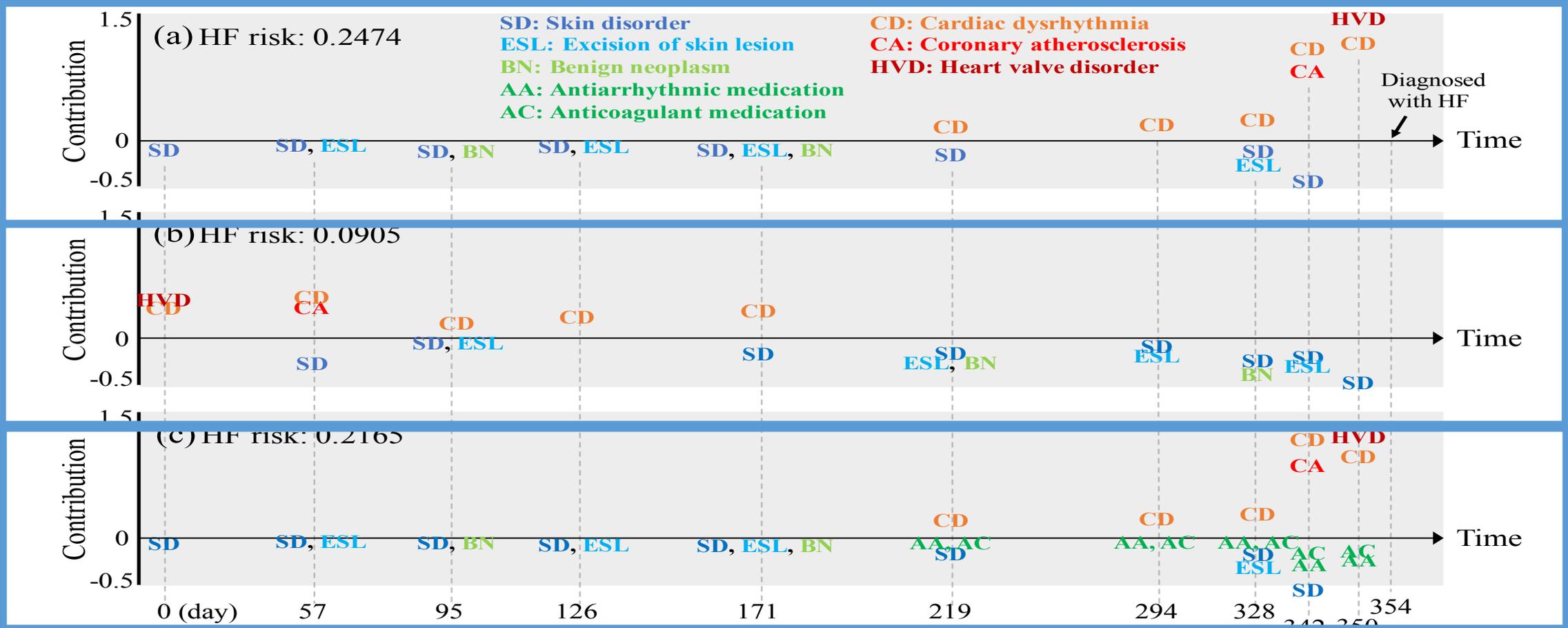
## Negative Log Likelihood on Test Set



## Classification AUC



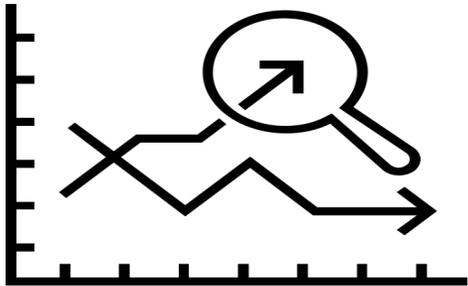
# RETAIN: Model Interpretability



# RETAIN: Summary



- Challenge: Deep learning models are often difficult to interpret



- RETAIN is a temporal attention model on electronic health records
  - Great predictive power
  - Good interpretation

# RAIM: Recurrent Attentive and Intensive Model of Multimodal Patient Monitoring Data

Yanbo Xu, Siddharth Biswal, Shriprasad Deshpande, Kevin Maher,  
and Jimeng Sun

**KDD' 18**

# RAIM: Motivation

**R**ecurrent  
**A**ttentive  
**I**ntensive  
**M**odel

1

Multi-channel high-density signal processing

2

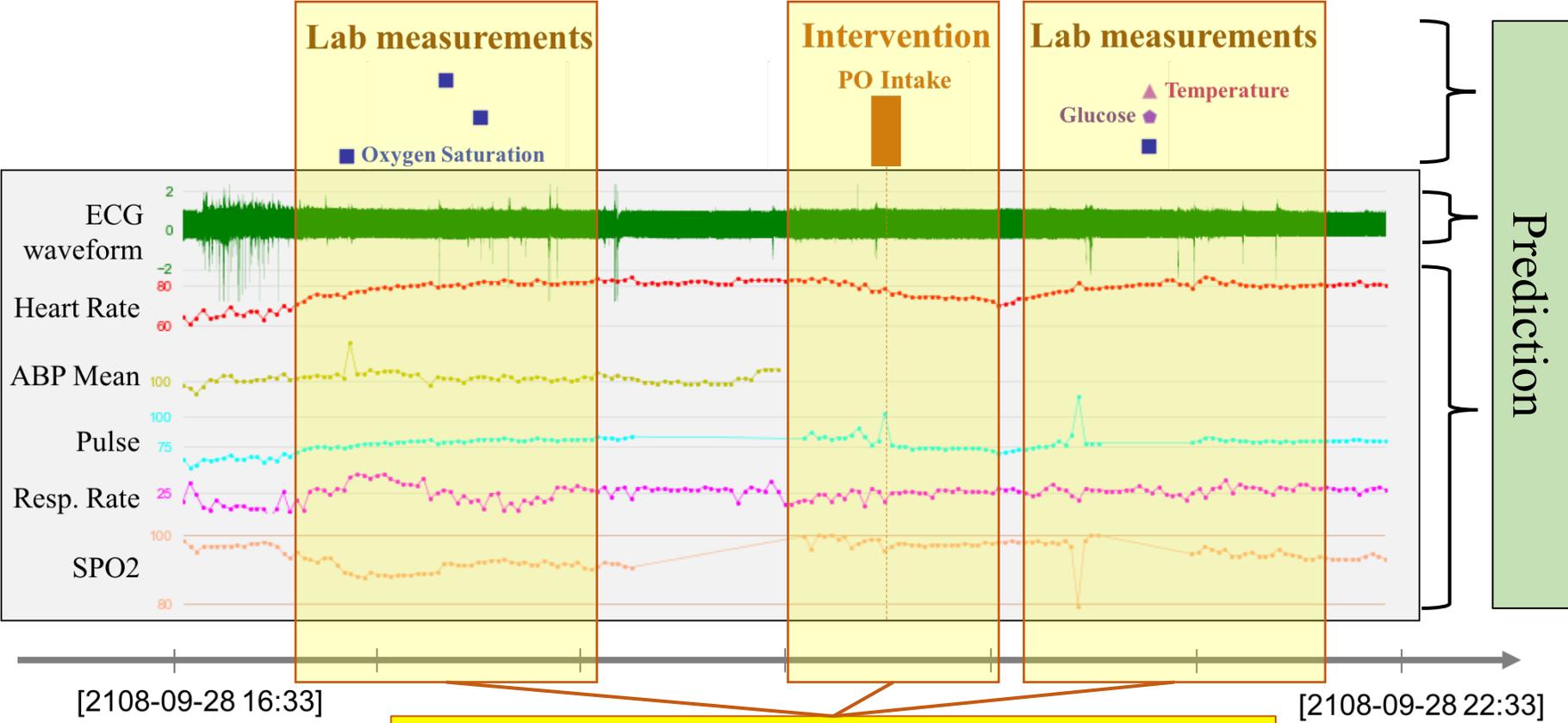
Multiple data modalities

3

Interpretability



# Challenges: interpretable predictions



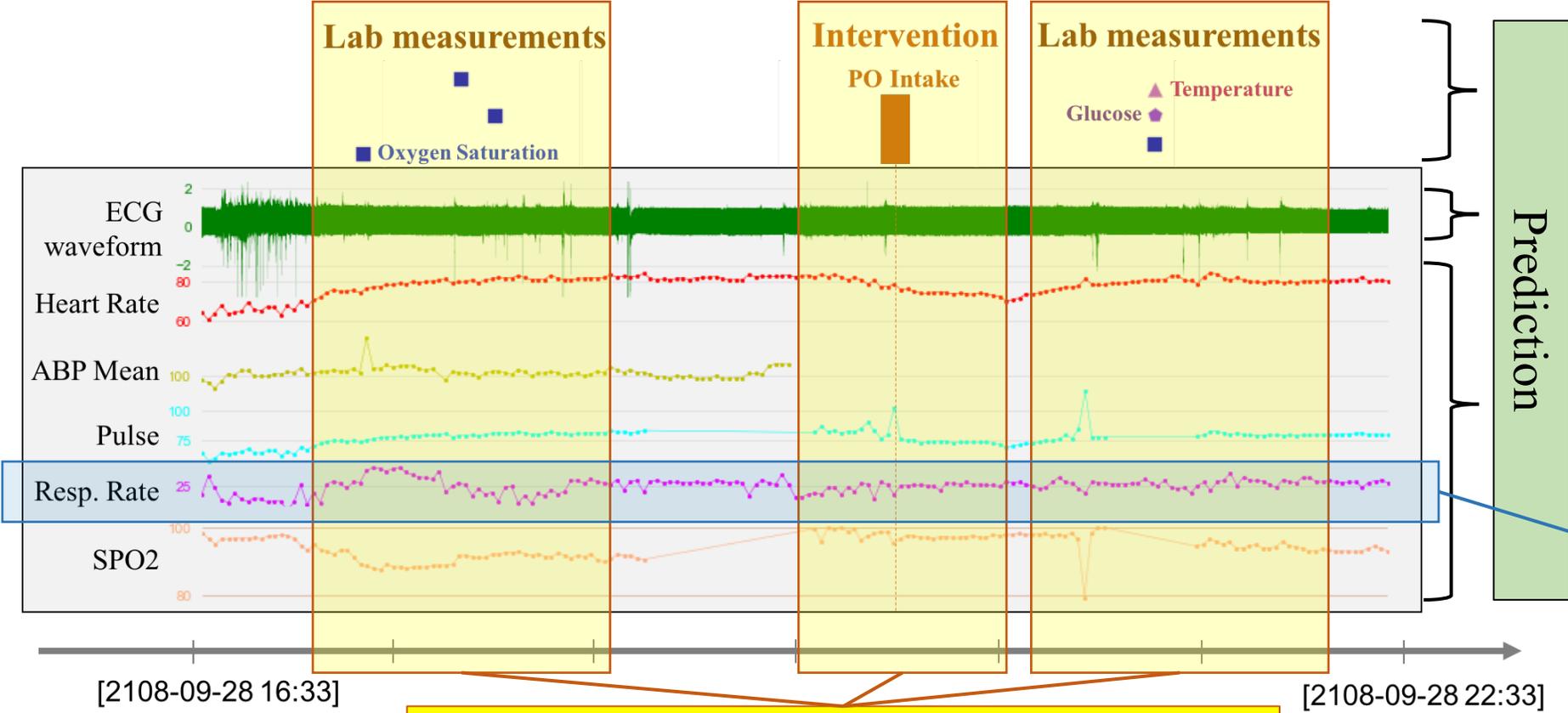
The predicted risk of decompensation is 0.84.



An example

More informative and probably influence most on the final prediction?

# Challenges: interpretable predictions



The predicted risk of decompensation is 0.84.

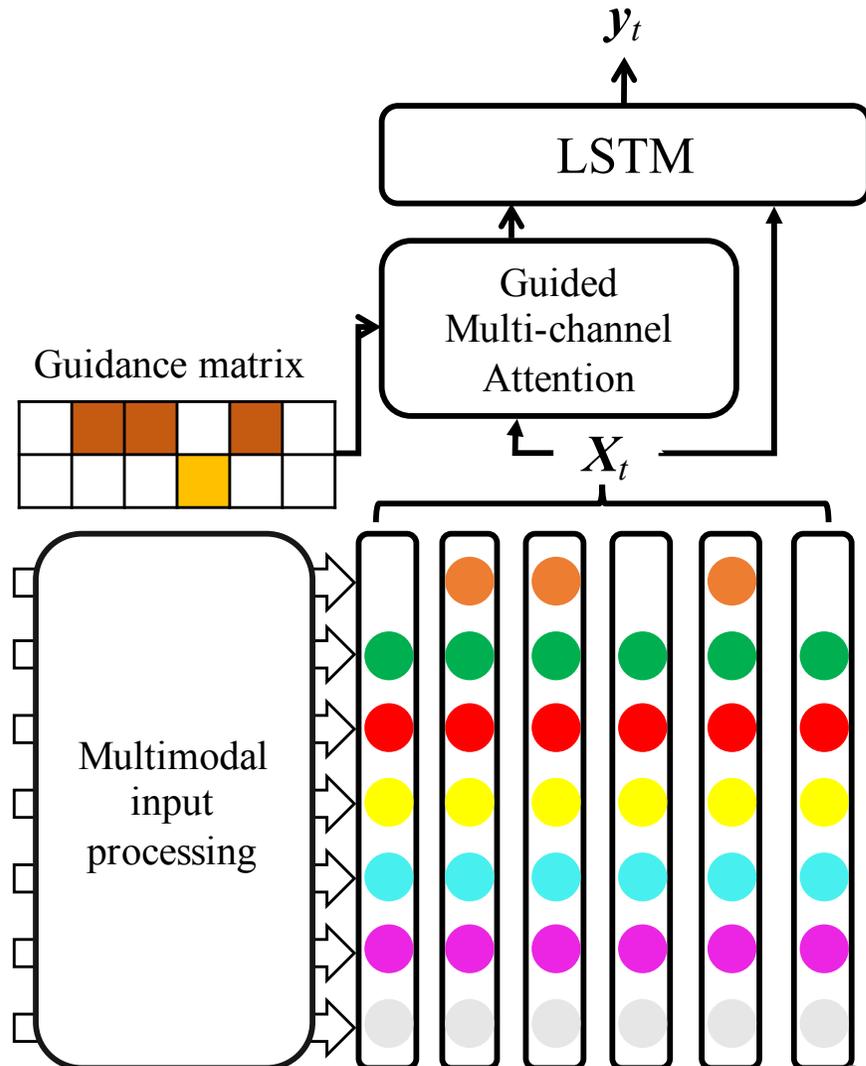
Which channel(s) are more important to attend for making the prediction?

More informative and probably influence most on the final prediction?



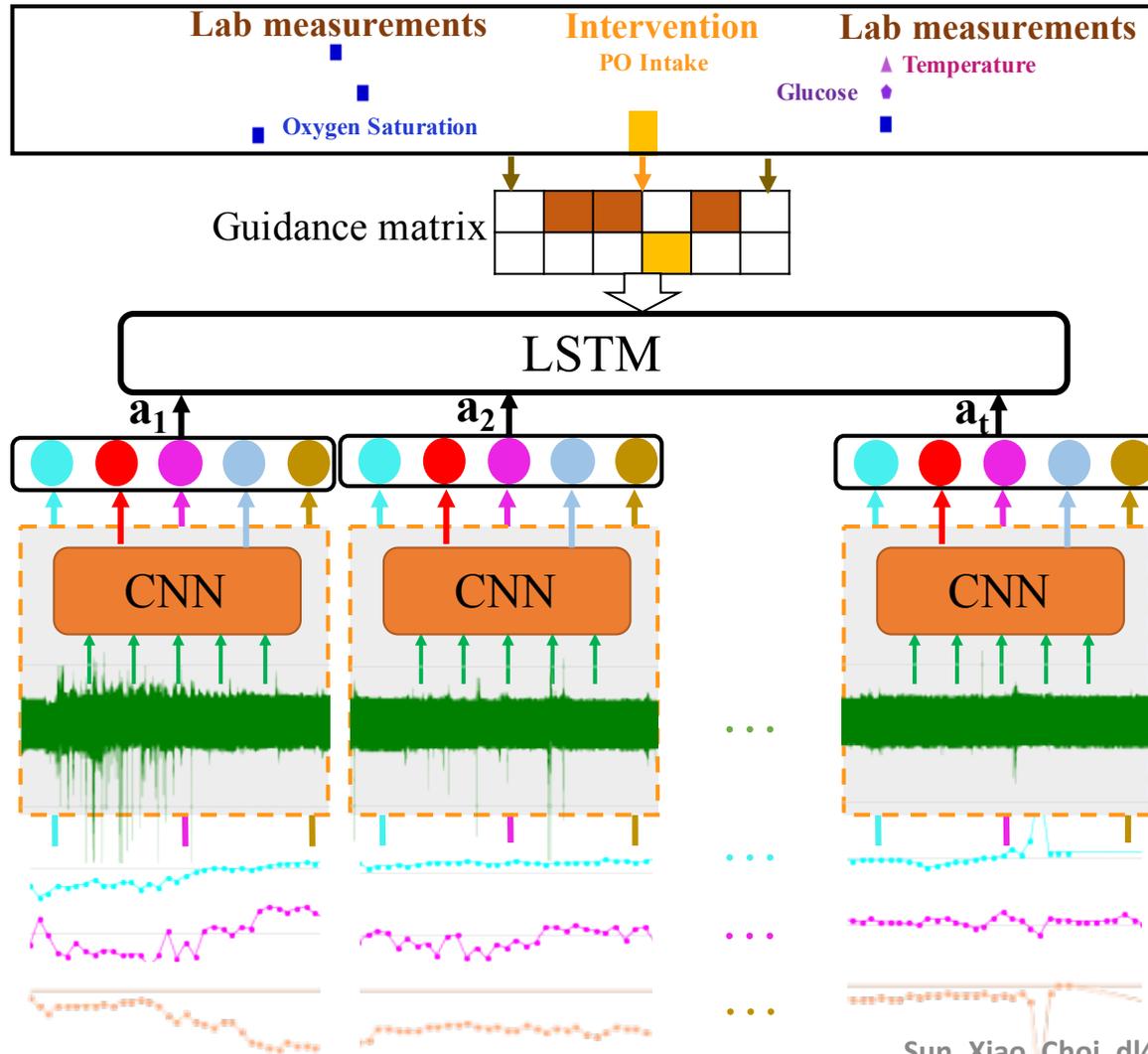
An example

# RAIM: A recurrent attentive and intensive model



- ✓ Multimodal input integrating
  - Irregular and regular discrete data
  - Continuous streaming data
- ✓ High-density signal processing
  - Recurrent deep neural networks
- ✓ Attention-based interpretable modeling
  - Multi-channel attention
  - Guided attention

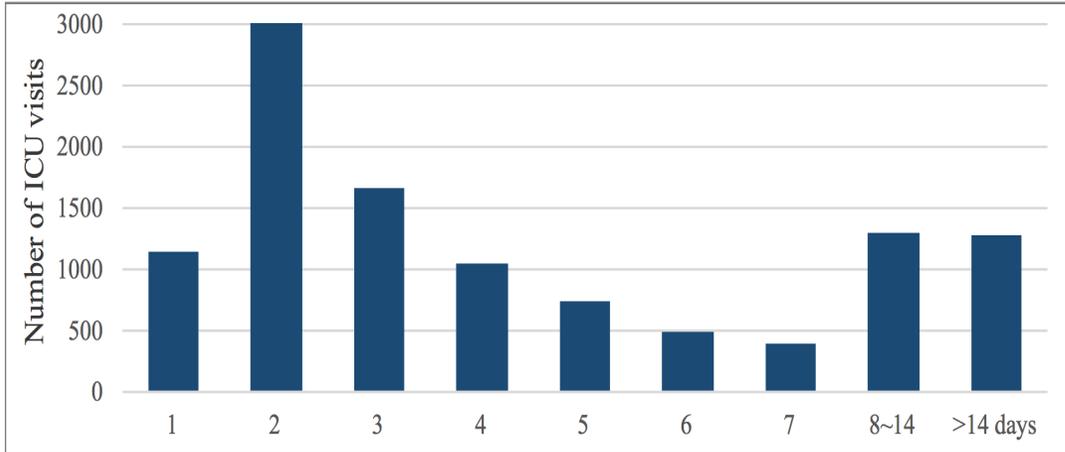
# RAIM: Multimodal input processing



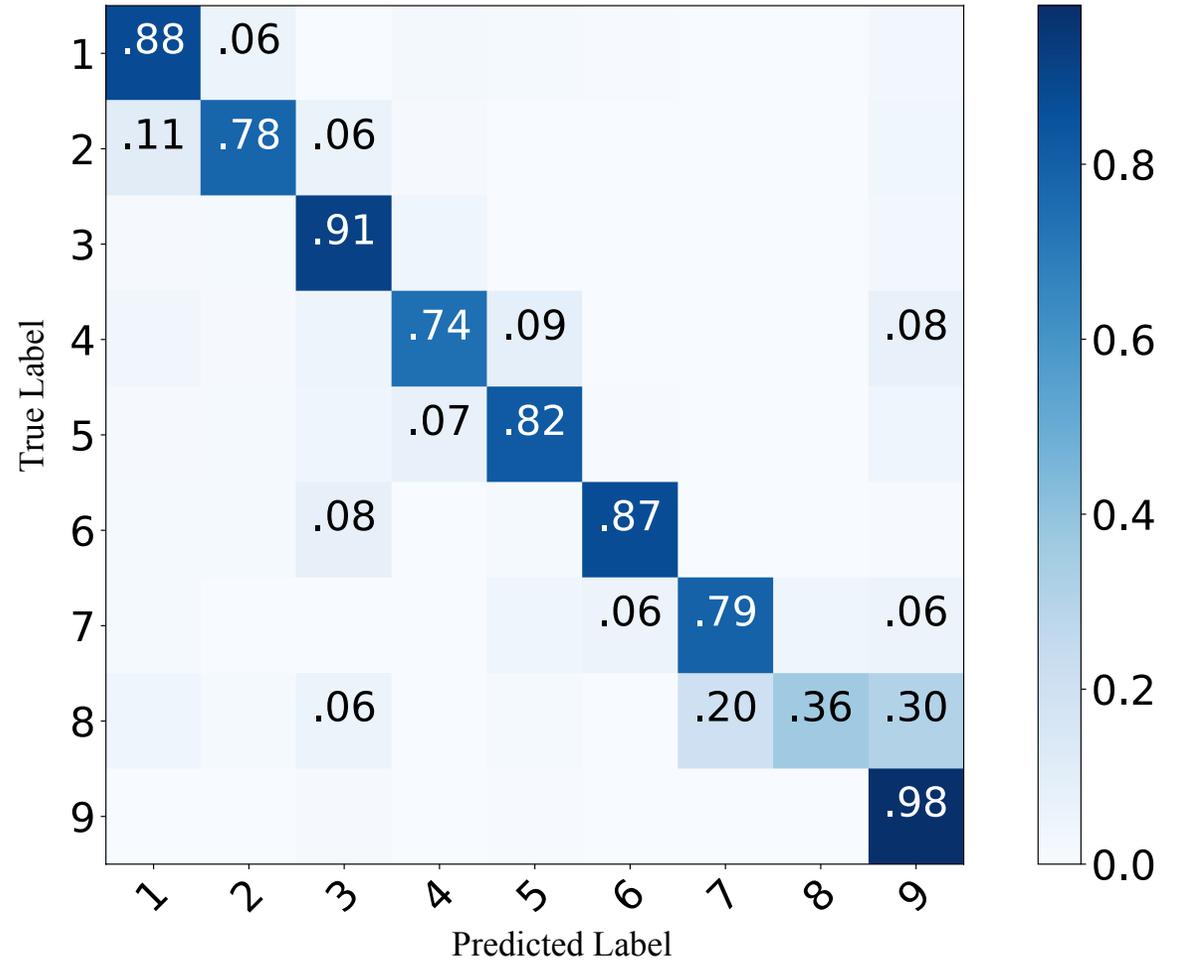
- Labs and medications are prioritized as a guidance matrix telling which time steps are more important to attend.
- RNN is applied for encoding long-term sequential behaviors.
- CNNs are applied for encoding short-term dense signals.
- Regularly recorded discrete values are unified into input vectors.

# RAIM: Predicting Length of Hospital Stay

(a) Histogram of length-of-stay



Number of cases



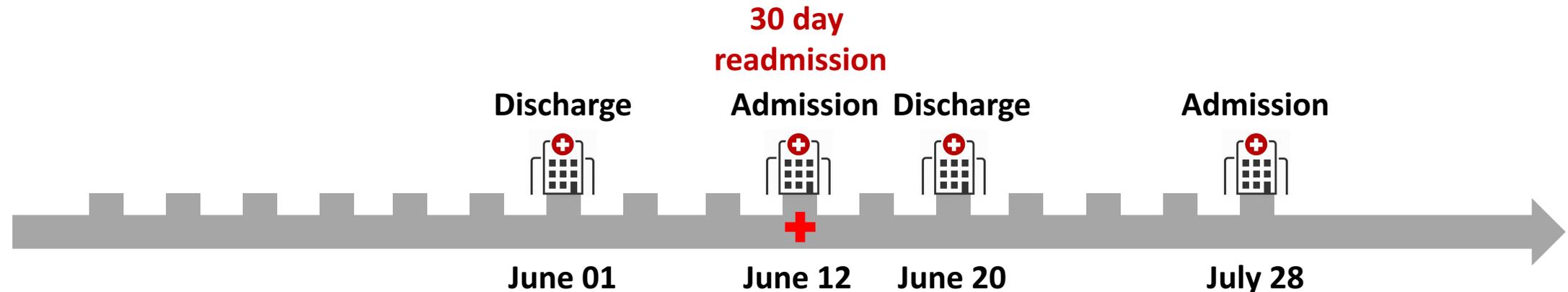
	Decompensation			Length of Stay	
	AUC-ROC	AUC-PR	Accuracy	Kappa	Accurac
CNN (ECG)	87.84%	21.56%	88.38%	0.7681	82.16%
CNN-RNN	87.45%	23.19%	88.25%	0.8027	85.34%
CNN-AttRNN	88.19%	25.81%	89.28%	0.8186	84.89%
RAIM-0	87.81%	25.56%	88.96%	0.8125	85.84%
RAIM-1	88.25%	25.61%	88.91%	0.8215	86.74%
RAIM-2	88.77%	26.85%	90.27%	0.8217	85.21%
RAIM-3	<b>90.18%</b>	<b>27.93%</b>	<b>90.89%</b>	<b>0.8291</b>	<b>86.82%</b>

# Readmission prediction via deep contextual embedding of clinical concepts

Cao Xiao , Tengfei Ma , Adji B. Dieng, David M. Blei, Fei Wang

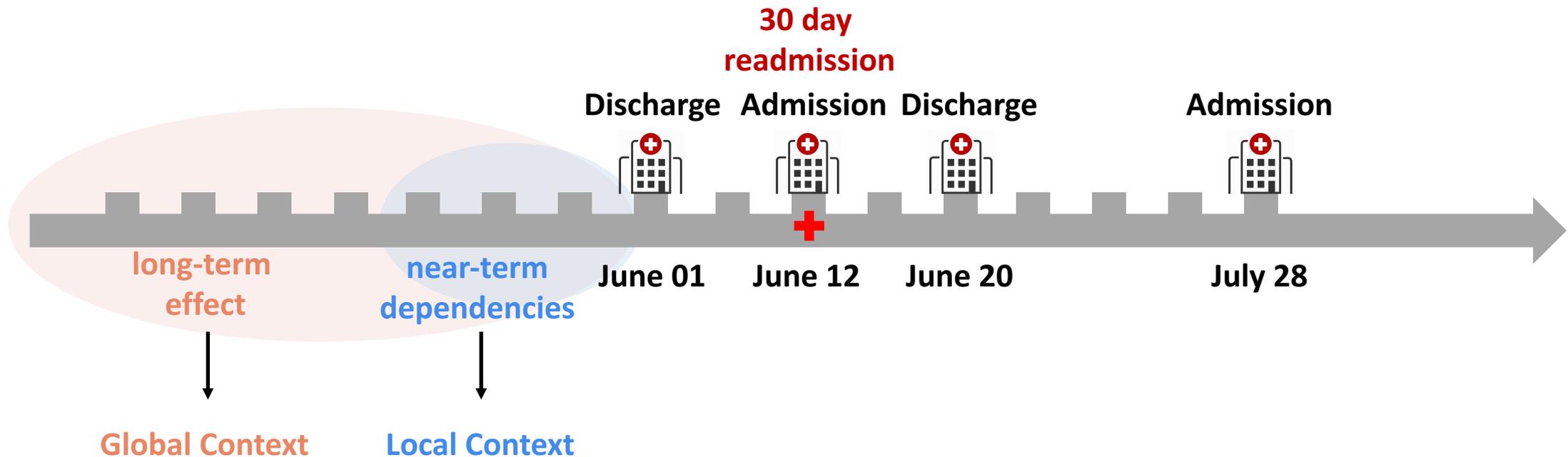
**PLOS One, 2018**

# Hospital Readmission



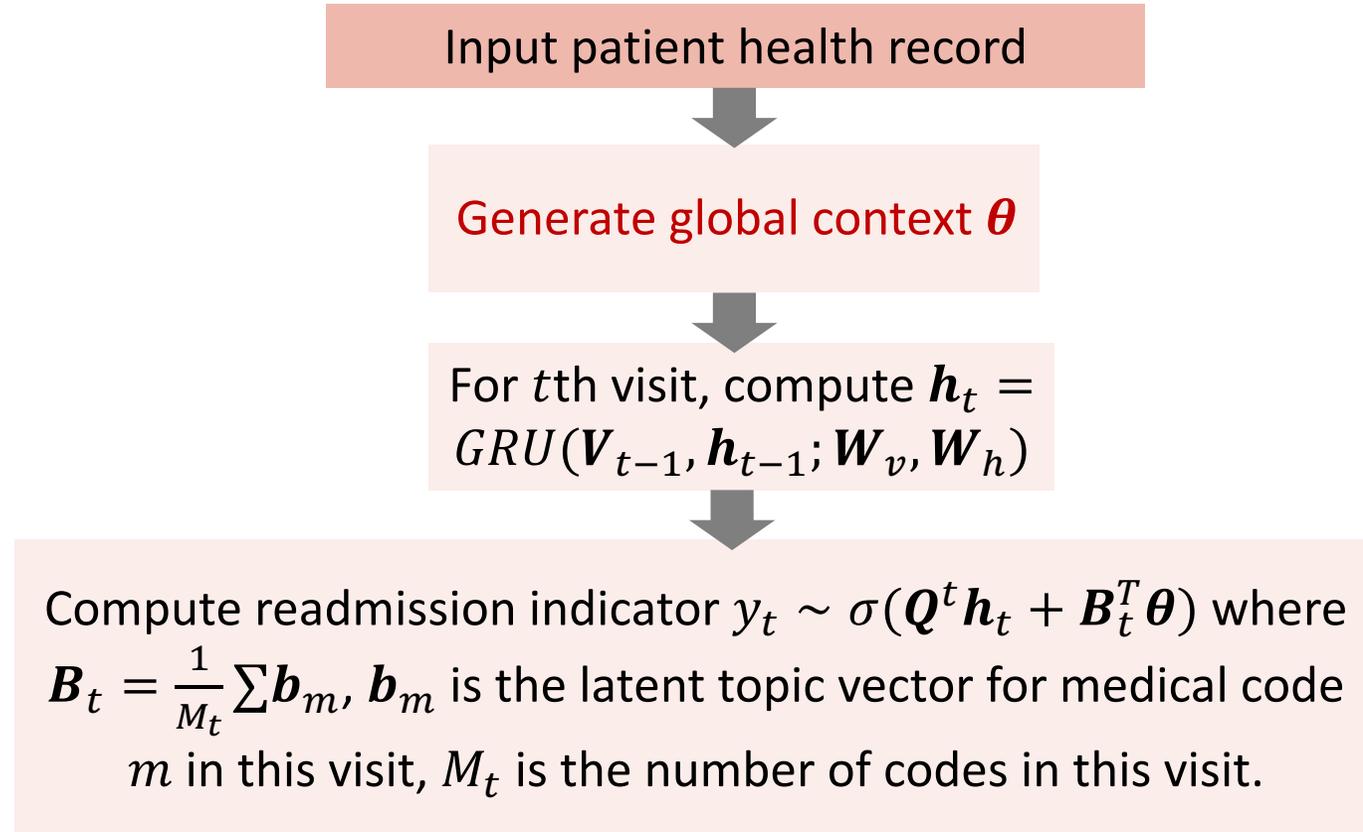
- In the U.S., involve 17.6% of hospital-admitted patients, accounted for \$17:9 billion Medicare spending per year, while 76% of them are potentially avoidable.
- Targeted follow-ups that focus on patients with high risks of readmissions could help reduce readmission rate.
- Readmission risks are **hard to predict** due to its **complex entanglements** with the patients' health conditions.

# CONTENT: Deep Contextual Embedding of Clinical Concepts



**CONTENT** is an end-to-end hybrid deep learning model structure that combines **topic modeling** and **Recurrent Neural Network (RNN)** to distill the complex knowledge hidden in those contexts.

# The CONTENT Model



# CONTENT: Model Inference

**Original  
Objective**

$$\log p(\mathbf{y}|\mathbf{h}, \Theta) = \log \int p(\mathbf{y}|\mathbf{h}, \Theta, \theta) p(\theta) d\theta.$$

**Approximate  
Objective**

$$\log p(\mathbf{y}|\mathbf{h}, \Theta) = D_{KL}(q(\theta) || p(\theta|\mathbf{y})) + ELBO.$$

**New  
Objective**

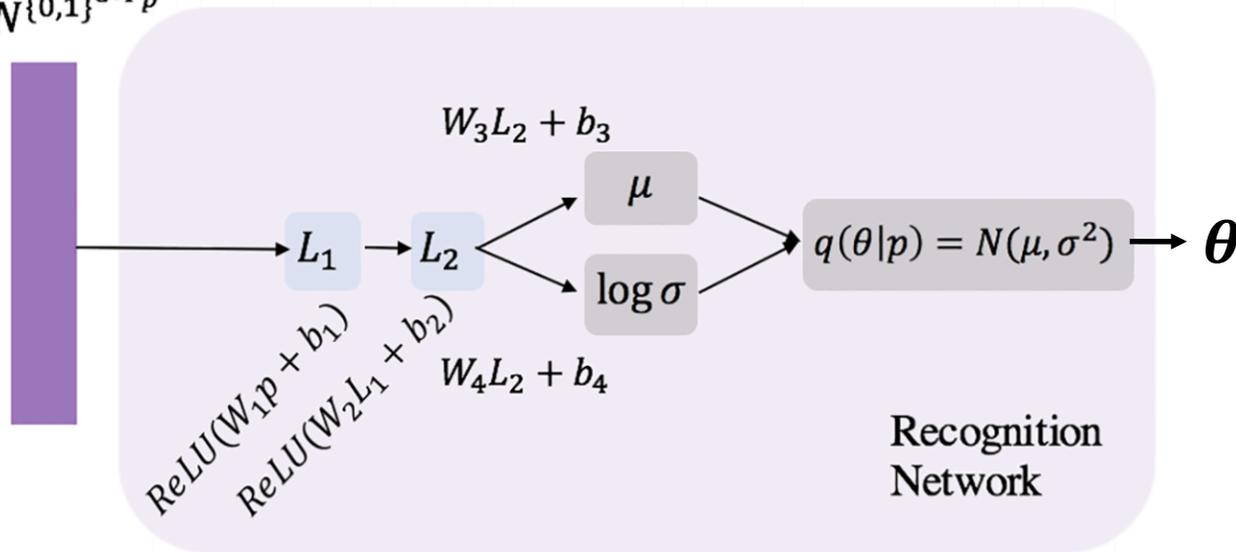
$$ELBO = E_{q(\theta)}[\log p(\mathbf{y}|\mathbf{h}, \theta, \Theta) + \log p(\theta) - \log q(\theta)] \leq \log p(\mathbf{y}|\mathbf{h}, \Theta)$$

**Inference Network**

# CONTENT: Inference Network $q(\theta)$

**Input:** Patient Representation Matrix

$$p \in N\{0,1\}^{C \times T_p}$$



**Generate Global Context Vector  $\theta$**

$$r_1 = \text{ReLU}(W_{r_1} p + b_{r_1})$$

$$r_2 = \text{ReLU}(W_{r_2} r_1 + b_{r_2})$$

$$\mu(p) = W_{\mu} r_2 + b_{\mu}$$

$$\log \sigma(p) = W_{\sigma} r_2 + b_{\sigma}$$

$$q(\theta|p) = N(\mu(p), \text{diag}(\sigma^2(p)))$$

$$\theta \sim q(\theta|p).$$

# CONTENT: Readmission Prediction

Dataset	Congestive Heart Failure
# patients	5, 393
# visits	455, 106
# events	1, 306, 685
Avg. # of visits per patient	84.4
Avg. # of events per patient	242.3
# of unique event codes	618

## Data split

- training: 4000 patients
- validation: 700 patients
- testing: 693 patients

Method	PR-AUC	ROC-AUC	ACC
Word2vec+LR	0.3445±0.0204	0.5360±0.0246	0.6828±0.0120
Med2vec+LR	0.3836±0.0149	0.5937±0.0120	0.6915±0.0095
GRU	0.3862±0.0136	0.5998±0.0124	0.6856±0.0082
GRU+Word2Vec	0.3430±0.0157	0.5616±0.0157	0.6731±0.0091
RETAIN	0.3720±0.0148	0.5707±0.0140	0.6814±0.0111
CONTENT	0.3894±0.0153	0.6103±0.0130	0.6934±0.0090

# CONTENT: Interpretability



Cluster 1



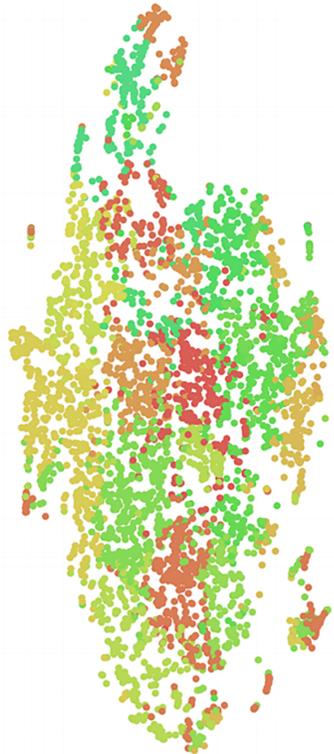
Cluster 2



Cluster 3



Cluster 4



Cluster 1: non-severe non-cardiac comorbidity group  
(average # readmission = 13.20)

Count	Name of Clinical Event
51	anemia
30	disorder of joint
28	disorder of back
26	osteoarthritis and allied disorders
19	symptoms involving respiratory system

Cluster 2: transplant surgery patient comorbidity group  
(average # readmission = 33.00)

Count	Name of Clinical Event
78	organ or tissue replaced by transplant
64	after-surgery care
63	acute renal failure
48	pneumonia
24	disorders of urethra and urinary tract

Cluster 3: cardiac-cancer comorbidity group  
(average # readmission = 21.09)

Count	Name of Clinical Event
965	diabetes mellitus
622	chronic airways obstruction
489	disorders of lipid metabolism
441	chronic ischemic heart disease
407	malignant neoplasm of female breast

Cluster 4: traumatic brain injury comorbidity group  
(average # readmission = 16.45)

Count	Name of Clinical Event
209	hypertensive heart disease
196	anemia
126	symptoms involving nervous and musculoskeletal system
125	intracranial injury
114	symptoms involving head and neck

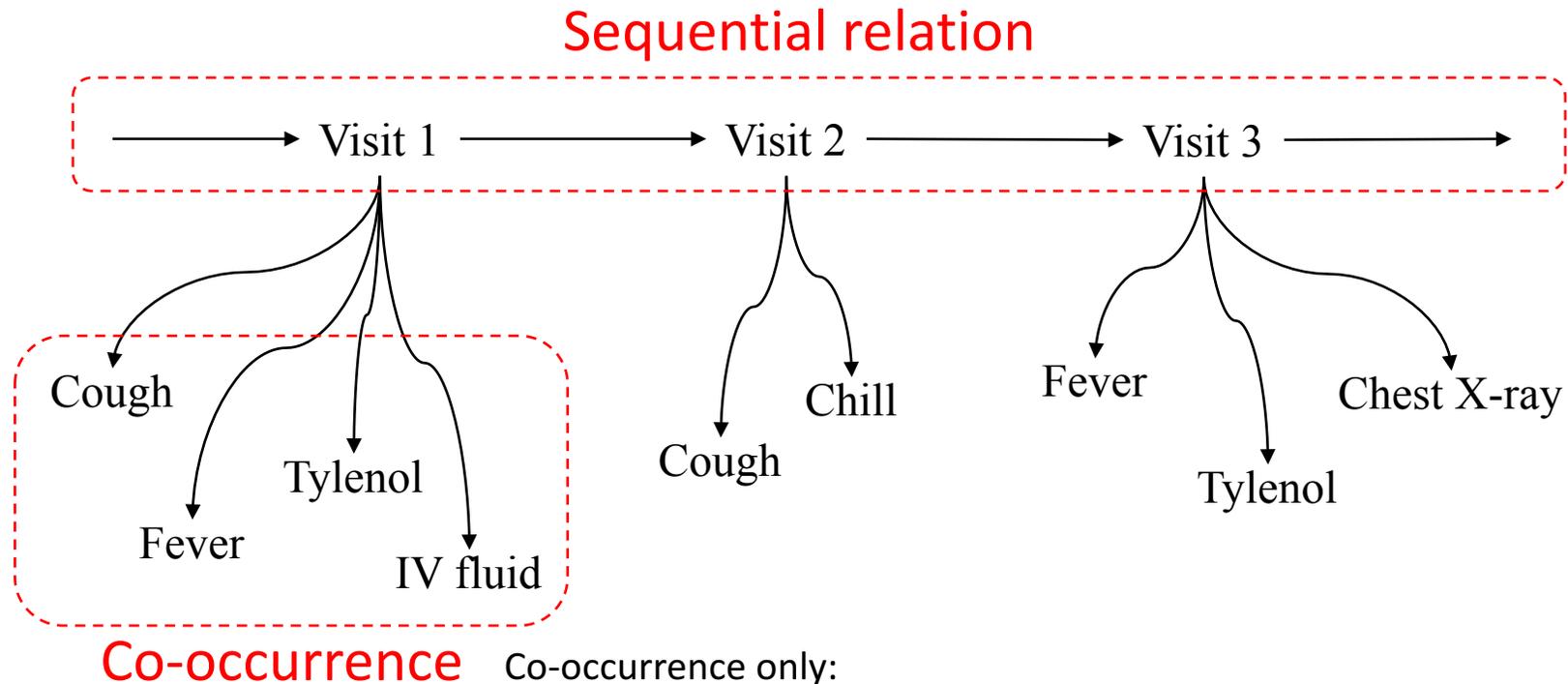
# Med2Vec: Multi-layer Representation Learning for Medical Concepts

Edward Choi, Mohammad T. Bahadori, Elizabeth Searles,  
Catherine Coffey, Jimeng Sun

**KDD'16**

# Med2Vec: Background

- Learn good representations of medical concepts
  - Diagnosis/medication/procedure codes
- Utilize 2-level structure of EHR



Co-occurrence only:

Choi, Youngduck, Chill Yi-I. Chiu, and David Sontag.

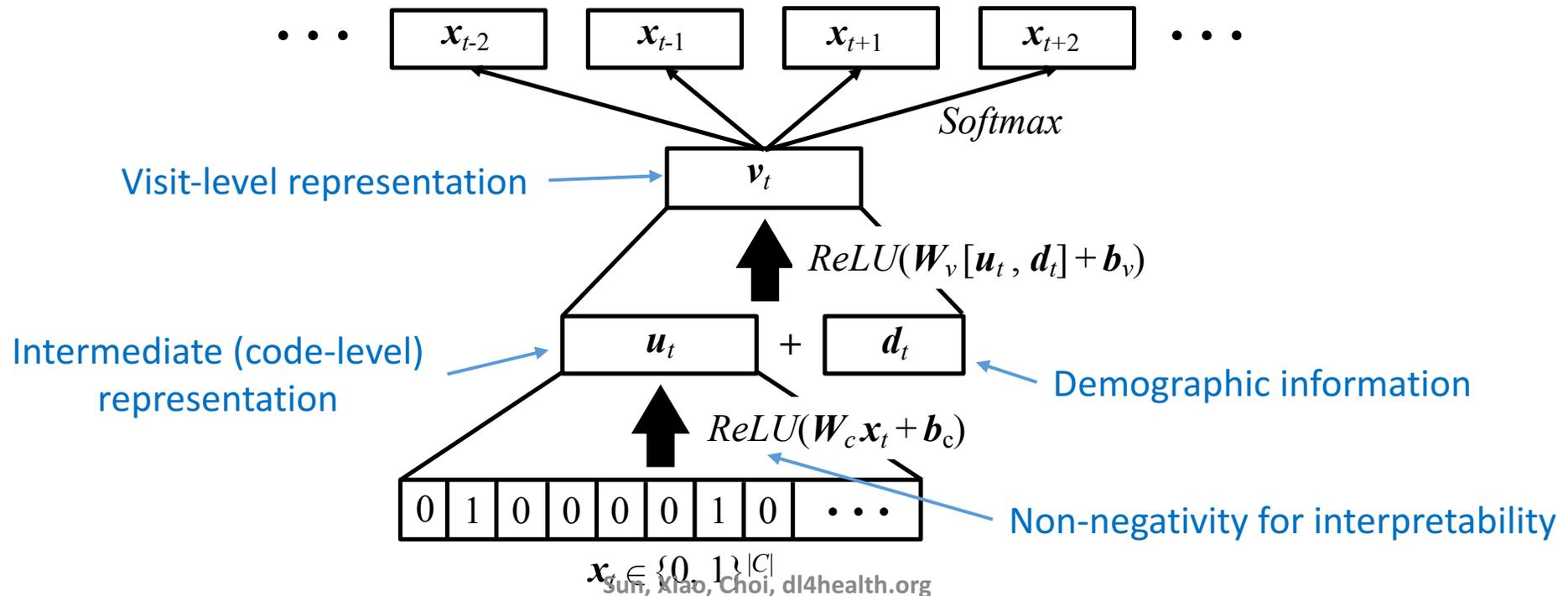
Sun, Xiao, Choi, dl4health.org

"Learning low-dimensional representations of medical concepts.", AMIA 2016

# Med2Vec: Model

- Model architecture

- Exploit two-layer structure of longitudinal EHR
  - Intra-visit codes provide co-occurrence information
  - Visit sequence provides sequential information



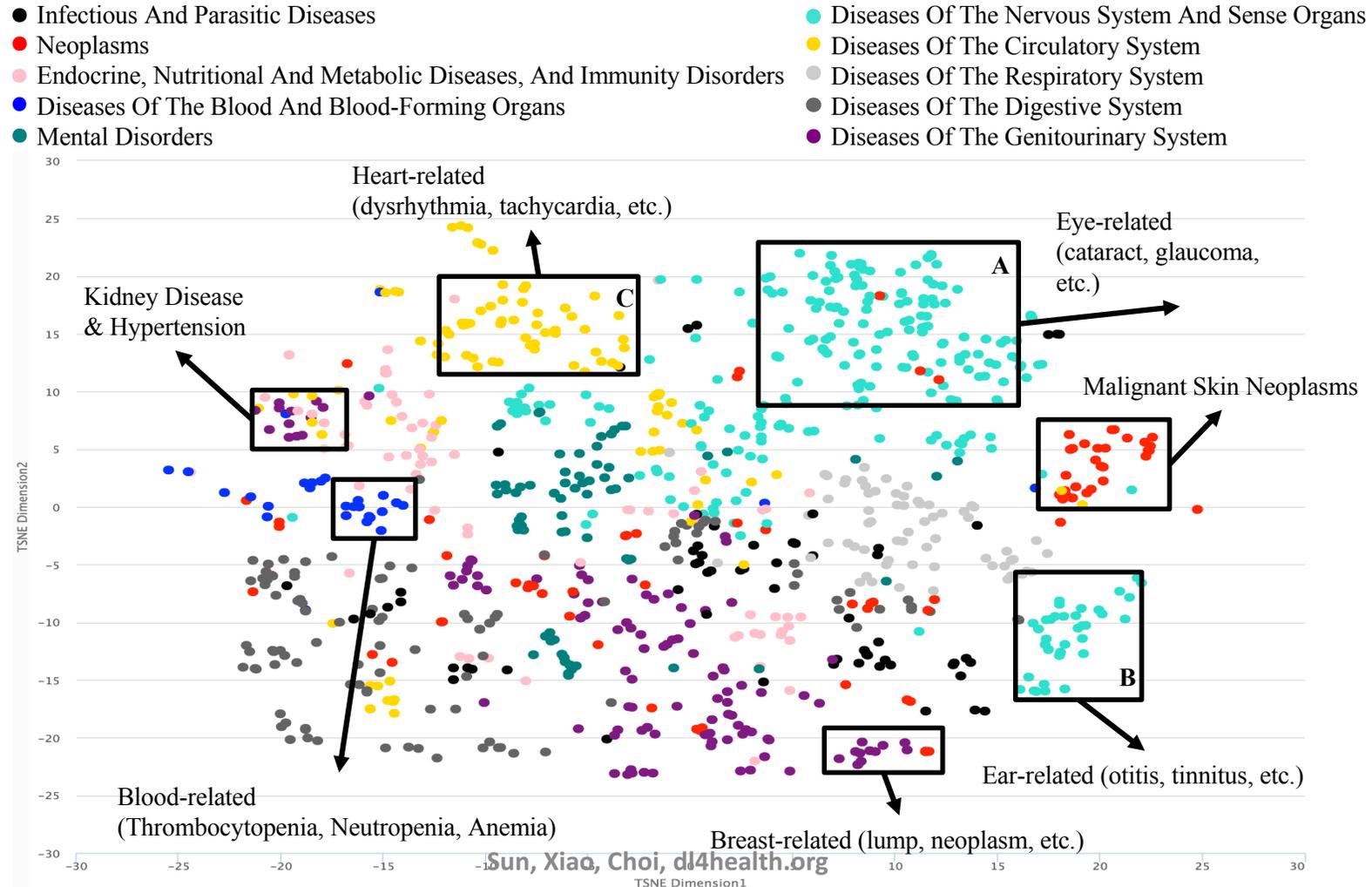
# Med2Vec: Data

- Data
  - Children's Healthcare of Atlanta

Dataset	CHOA
# of patients	550,339
# of visits	3,359,240
Avg. # of visits per patient	6.1
# of unique medical codes	28,840
- # of unique diagnosis codes	10,414
- # of unique medication codes	12,892
- # of unique procedure codes	5,534
Avg. # of codes per visit	7.88
Max # of codes per visit	440
(95%, 99%) percentile # of codes per visit	(22, 53)

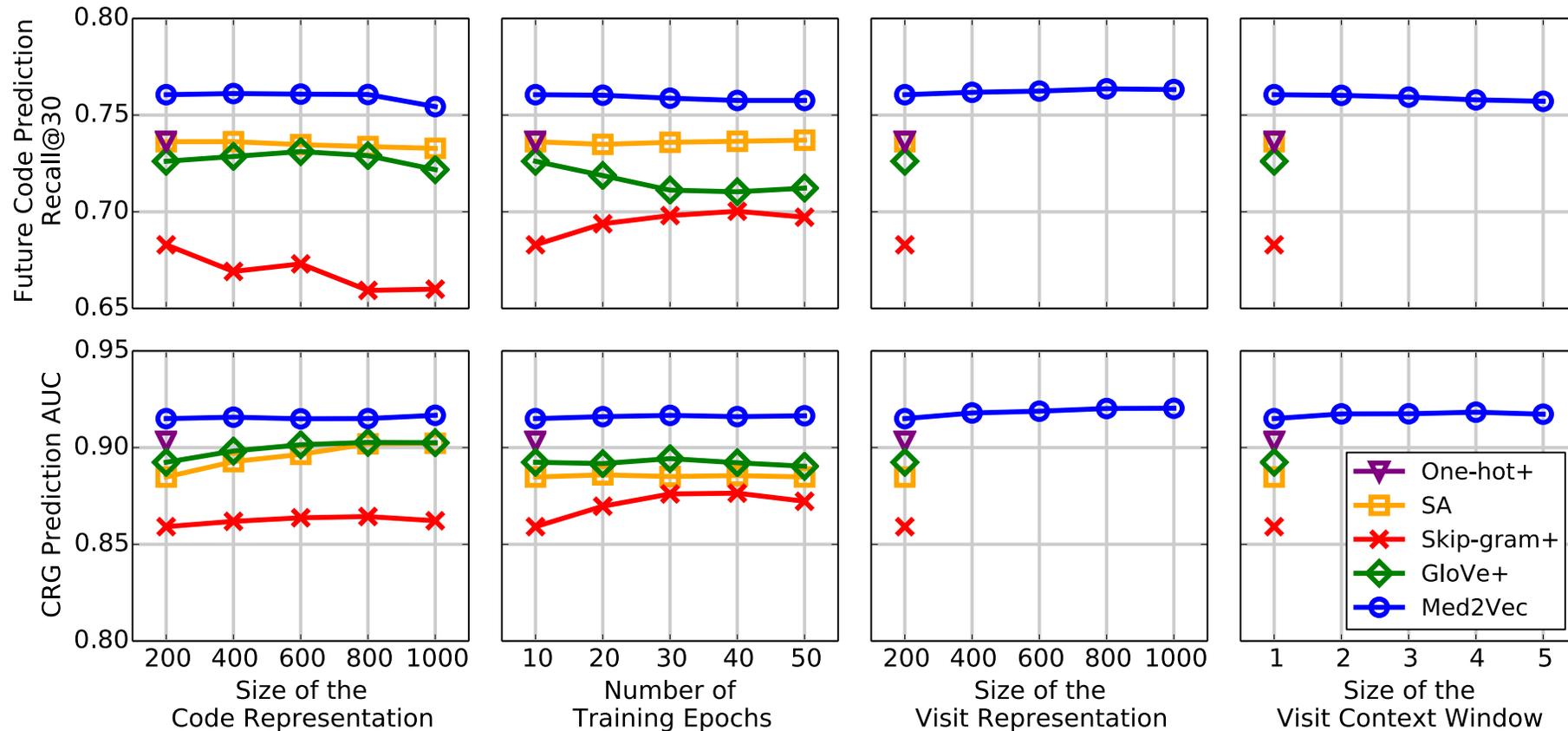
# Med2Vec: Result

- Visualizing the learned representations



# Med2Vec: Prediction

- Using the learned representations for prediction



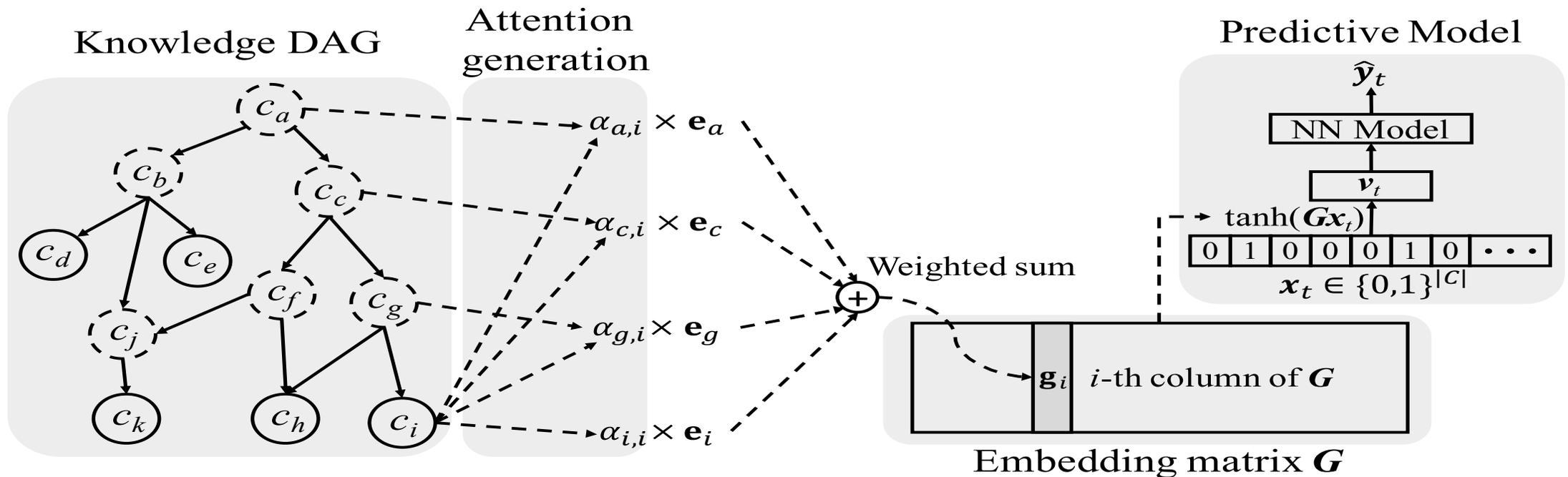
# GRAM: Graph-based Attention Model for Healthcare Representation Learning

Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F. Stewart, Jimeng Sun

**KDD' 17**

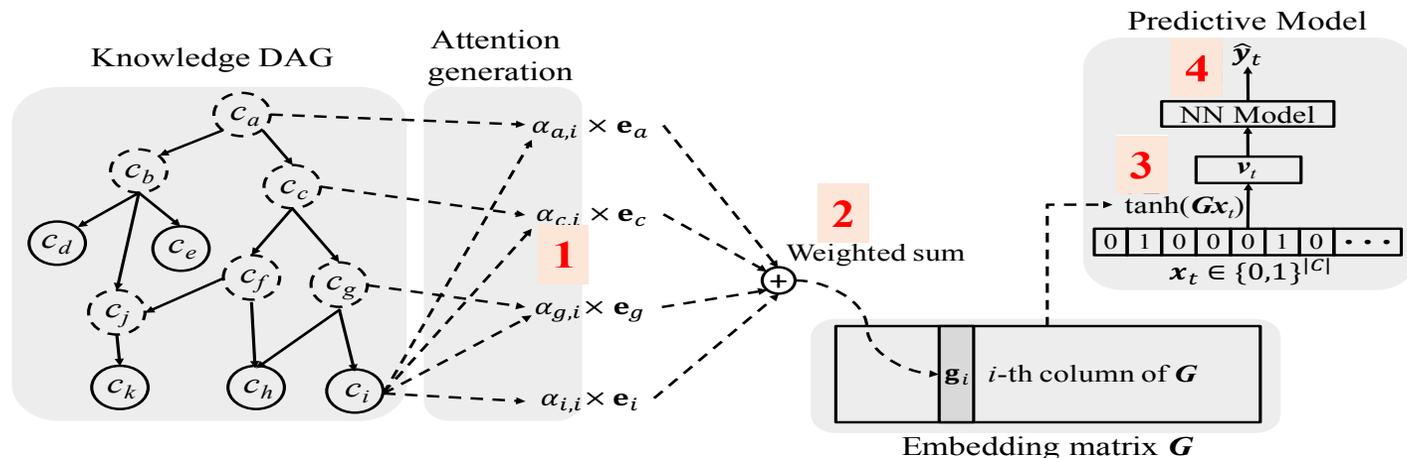
# GRAM: Learn representations of medical codes leveraging medical ontologies

- Method: Generate a medical code representation vector by combining the representation vectors of its ancestors using the attention mechanism



Model structure of GRAM

# GRAM: Algorithm

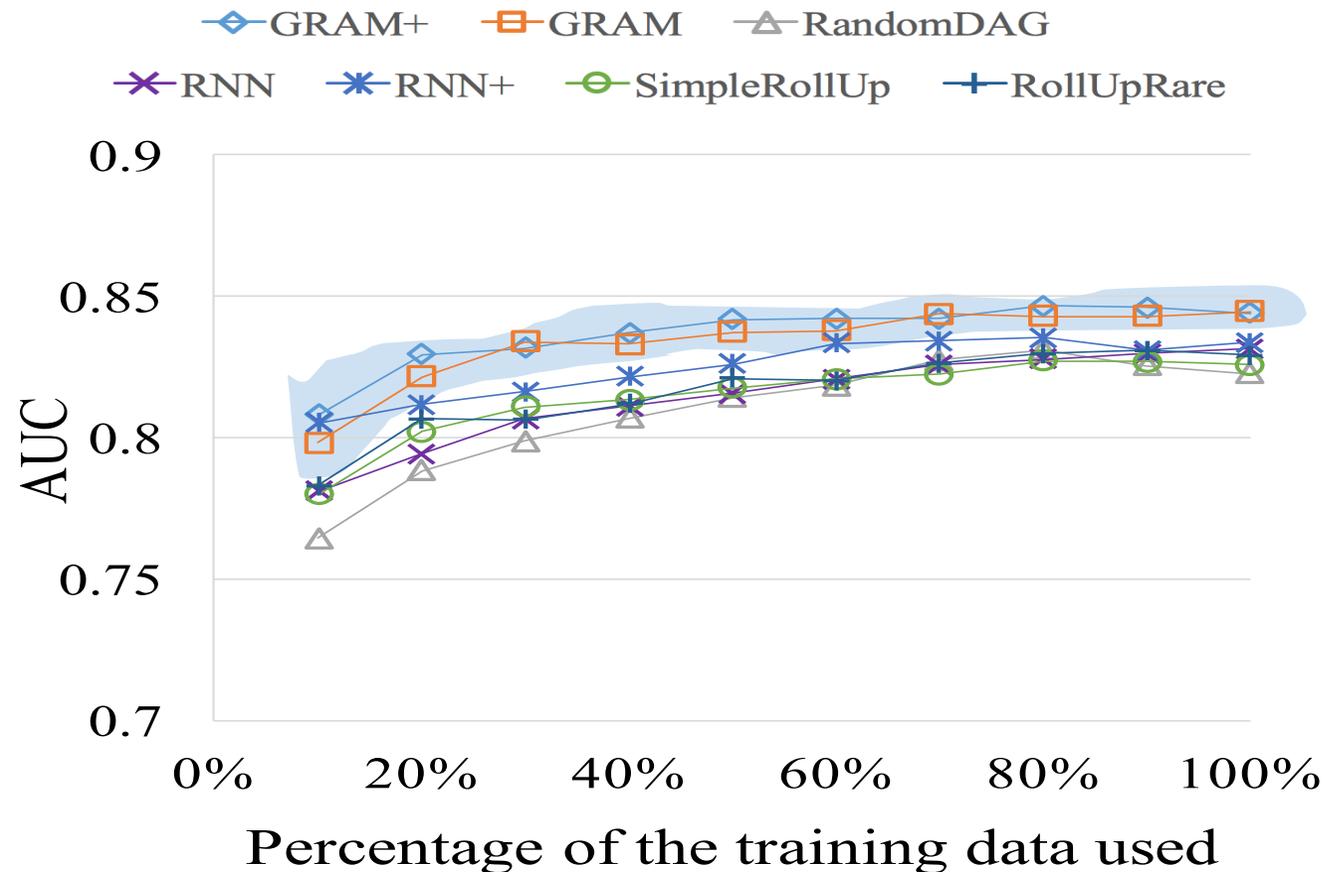


<b>1</b>	$\alpha_{ij} = \frac{\exp(f(\mathbf{e}_i, \mathbf{e}_j))}{\sum_{k \in \mathcal{A}(i)} \exp(f(\mathbf{e}_i, \mathbf{e}_k))}$ where $f(\mathbf{e}_i, \mathbf{e}_j) = \mathbf{u}_a^\top \tanh(\mathbf{W}_a \begin{bmatrix} \mathbf{e}_i \\ \mathbf{e}_j \end{bmatrix} + \mathbf{b}_a)$
	Attention weights are generated for all pairs of basic embeddings and its ancestors.
<b>2</b>	$\mathbf{g}_i = \sum_{j \in \mathcal{A}(i)} \alpha_{ij} \mathbf{e}_j,$
	Final representation is the weighted sum of attention weights and basic embeddings.
<b>3</b>	$\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_t = \tanh(\mathbf{G}[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t])$
	Sequence of visit representations are obtained using the Embedding matrix $G$ .
<b>4</b>	$\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_t = \text{RNN}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_t, \theta_r),$ $\hat{\mathbf{y}}_t = \hat{\mathbf{x}}_{t+1} = \text{Softmax}(\mathbf{W}\mathbf{h}_t + \mathbf{b}),$
	Performing sequential diagnoses prediction, outcomes are generated by RNN and Softmax.

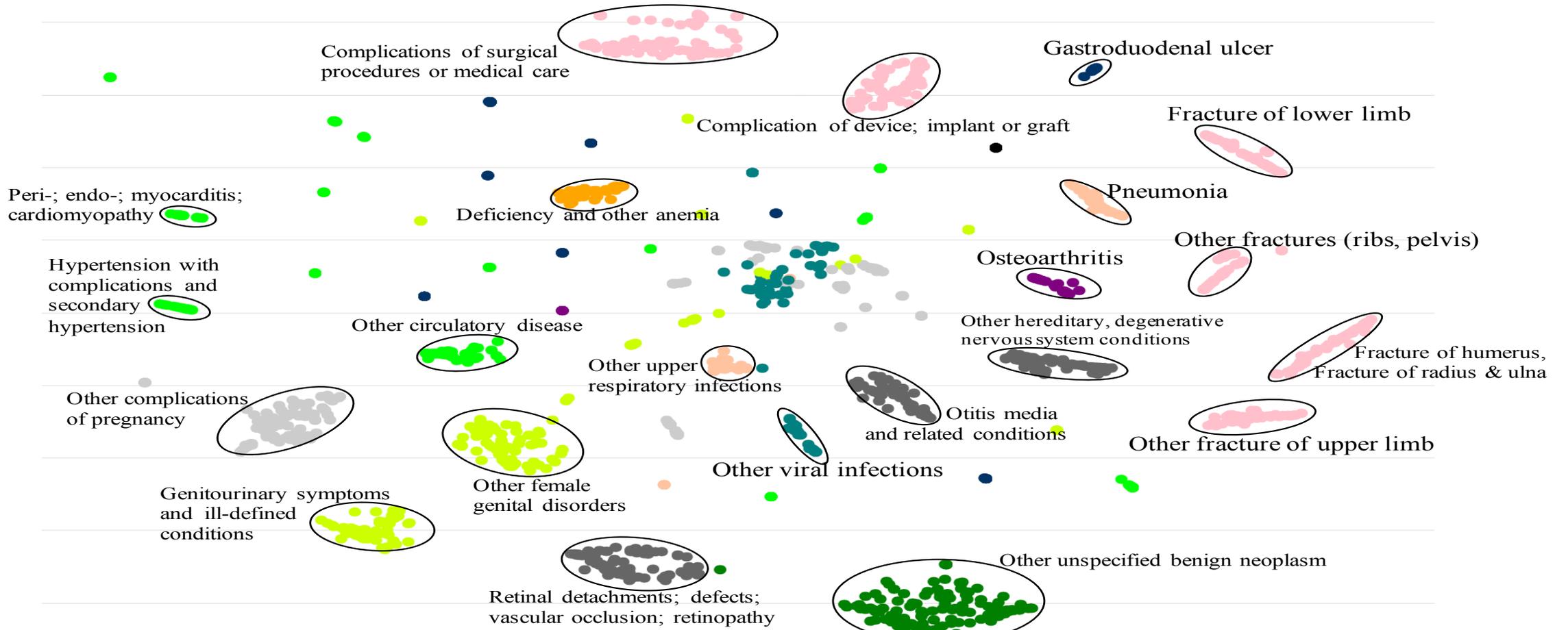
# GRAM provide accurate prediction

GRAM shows better predictive performance under data constraints

HF prediction using varying sizes of training data

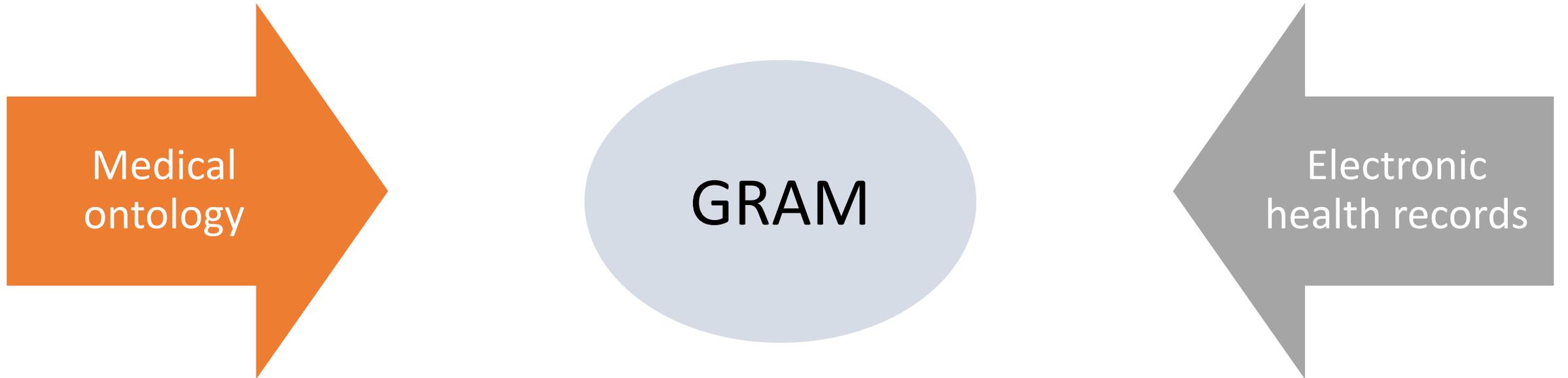


# GRAM learns representations well aligned with knowledge ontology



Scatterplot of GRAM representations

## GRAM: Summary



- Robust representation against **data insufficiency**
- Interpretable: Well aligned with medical knowledge

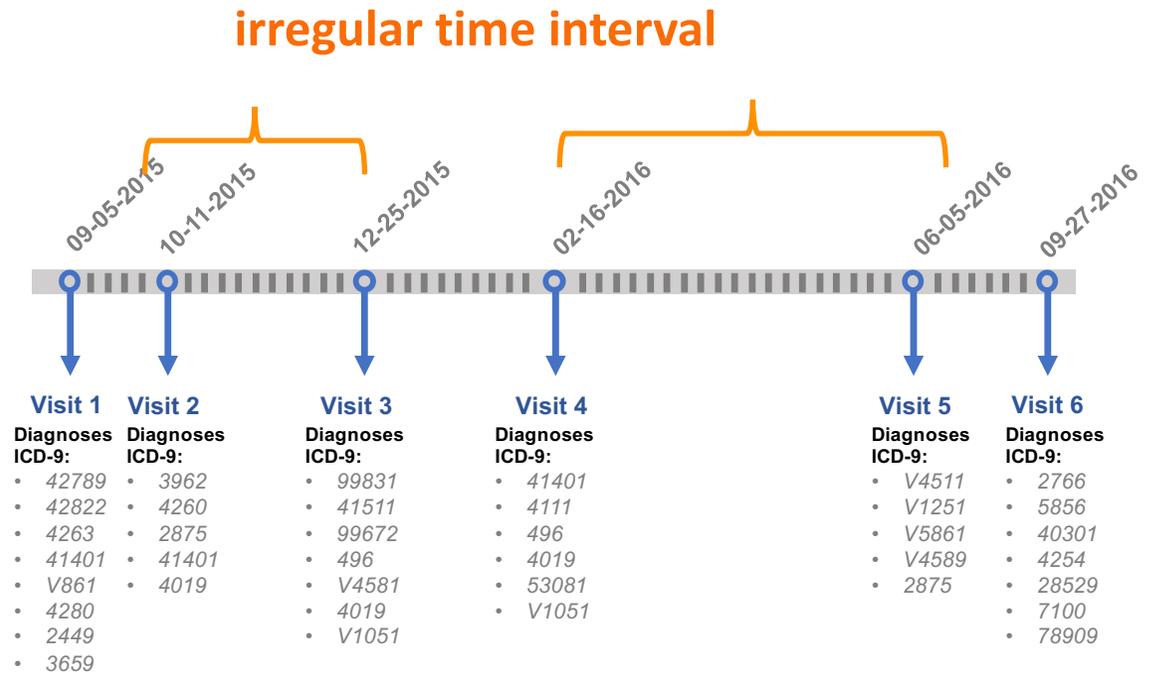
# T-LSTM: Patient Subtyping via Time-Aware LSTM Networks

Inci M. Baytas, Cao Xiao, Xi Zhang, Fei Wang, Anil K. Jain, Jiayu Zhou

**KDD' 17**

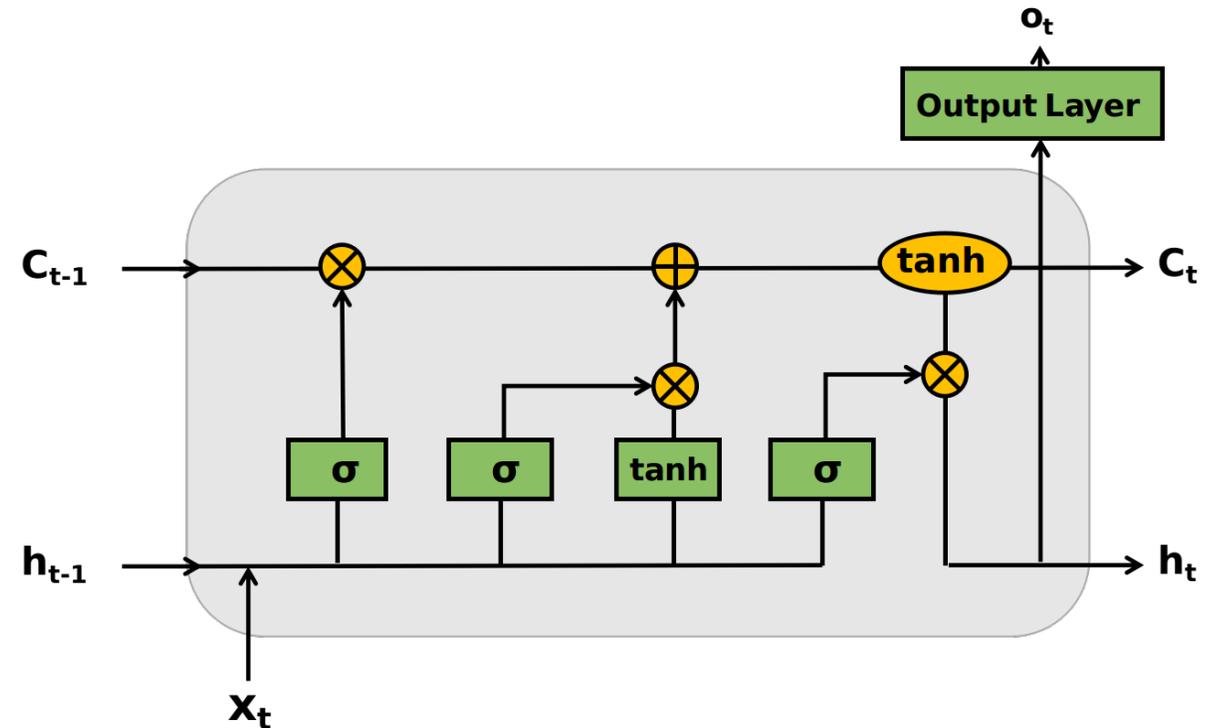
# T-LSTM: Background

- **Patient subtyping** seeks patient groups with similar disease progression pathways based on longitudinal EHR.
- **Elapsed time** has a significance in clinical decision making.
- It is crucial to **capture the relationships and the dependencies between clinical events under time interval irregularities.**

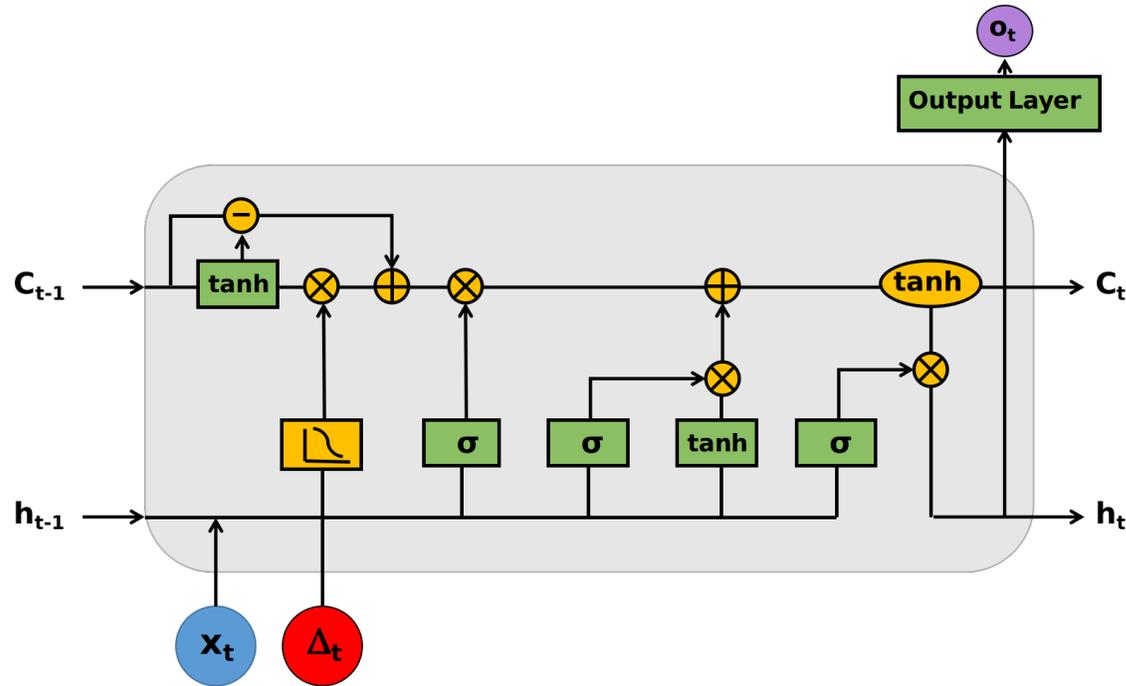


# T-LSTM: Limitation of LSTM

- The Long-Short Term Memory (LSTM) is an appealing model for capturing long-range dependency in EHR data.
- However, LSTM assumes the elapsed time is uniform throughout the sequence.



# The Time-Aware LSTM (T-LSTM) Unit



**Input Elapsed**  
**Records Time**

## Subspace Decomposition

$$C_{t-1}^S = \tanh(W_d C_{t-1} + b_d) \quad (\text{Short-term memory})$$

$$\hat{C}_{t-1}^S = C_{t-1}^S * g(\Delta_t) \quad (\text{Discounted short-term memory})$$

$$C_{t-1}^T = C_{t-1} - C_{t-1}^S \quad (\text{Long-term memory})$$

$$C_{t-1}^* = C_{t-1}^T + \hat{C}_{t-1}^S \quad (\text{Adjusted previous memory})$$

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (\text{Forget gate})$$

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (\text{Input gate})$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (\text{Output gate})$$

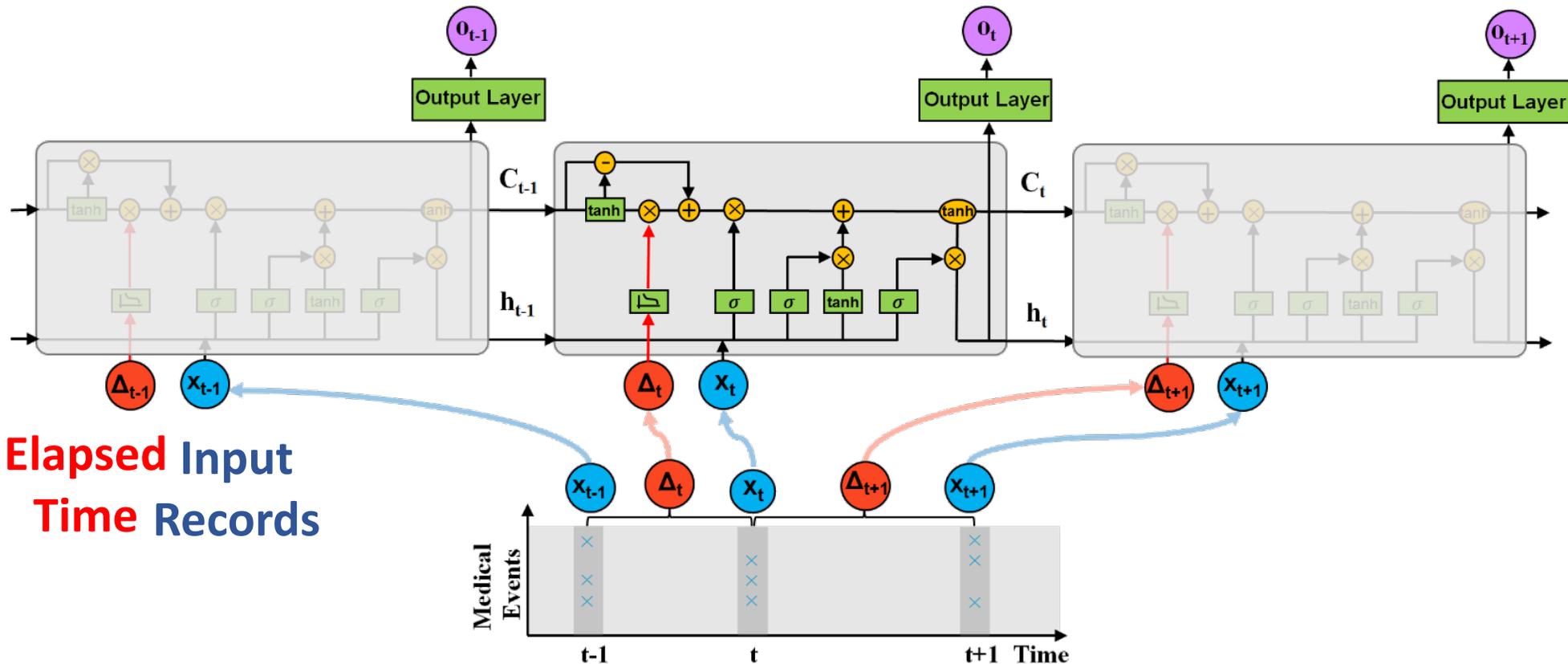
$$\tilde{C} = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (\text{Candidate memory})$$

$$C_t = f_t * C_{t-1}^* + i_t * \tilde{C} \quad (\text{Current memory})$$

$$h_t = o_t * \tanh(C_t), \quad (\text{Current hidden state})$$

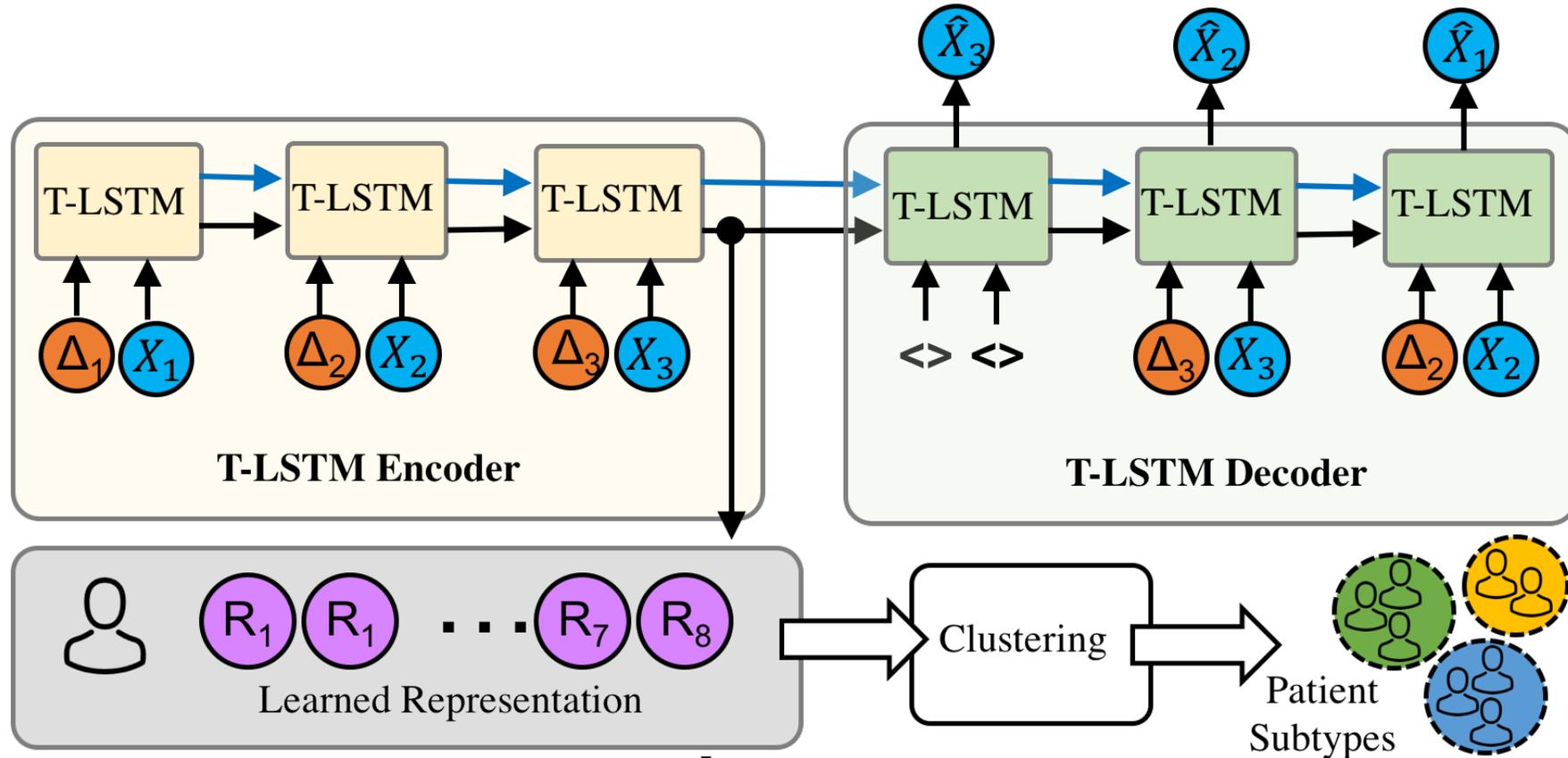
## Standard LSTM

# Time-Aware LSTM (T-LSTM)



T-LSTM decomposes the previous memory into long and short term components and utilizes the elapsed time ( $\Delta t$ ) to discount the short term effects.

# T-LSTM Autoencoder



$$E_r = \sum_{i=1}^L \left\| x_i - \hat{x}_i \right\|_2^2$$

# Time-Aware LSTM: Supervised Task

- Binary Classification of diabetes mellitus
- Clinical codes from 6730 patients
- Data are EMRbot
- <http://www.emrbots.org/>

Methods	Avg. Test AUC	Stdev.
T-LSTM	<b>0.91</b>	0.01
MF1-LSTM	0.87	0.02
MF2-LSTM	0.82	0.09
LSTM	0.85	0.02
LR	0.56	0.01

# Time-Aware LSTM: Parkinson's Disease Patient Subtyping

- Features include:
  - Demographics
  - Motor severity measures such as Unified Parkinson's Disease Rating Scale, Hoehn and Yahr staging,
  - Non-motor manifestations: depression, anxiety, cognitive status, sleep disorders,
  - Imaging assessment such as DaTScan.

Parkinson's Progression Markers Initiative (PPMI) data, 654 patients

Elapsed time interval: [1,26] months

Average sequence length: 25

Input feature dimension: 319

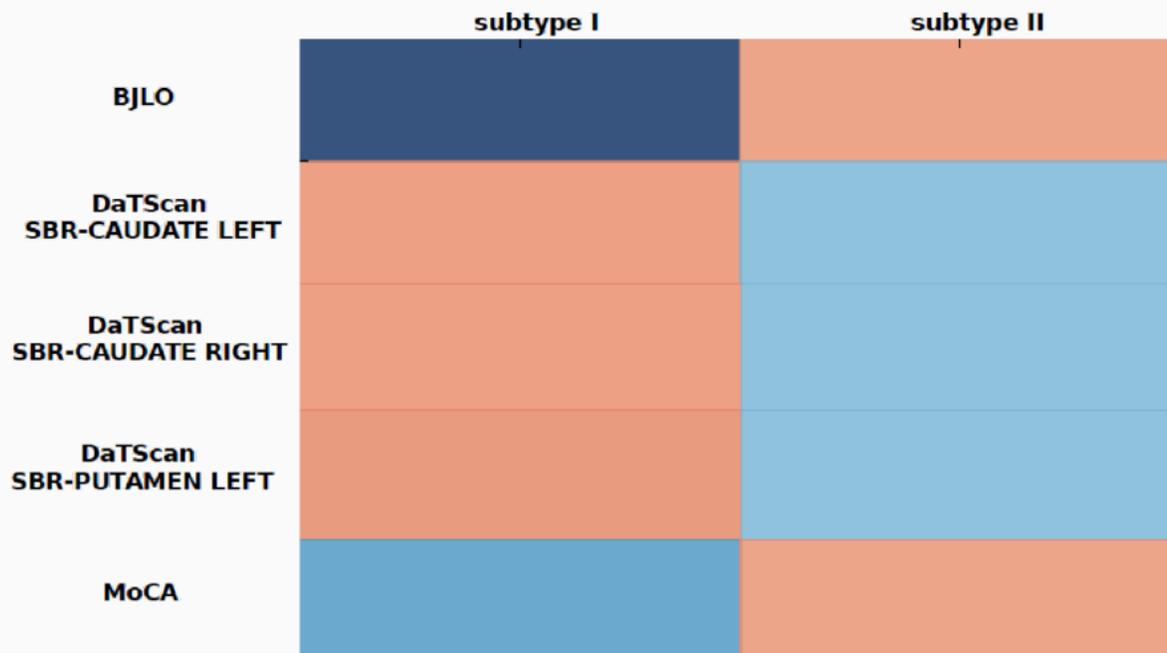
Target dimension: 82

# Time-Aware LSTM: Parkinson's Disease Patient Subtyping

Feature	P-Value	Cluster1 Mean	Cluster2 Mean
<b>T-LSTM</b>			
BJLO	$9.51 \times 10^{-8}$	16.5	24.7
MoCA	0.001	40.0	41.2
DaTScan1	0.042	2.29	2.07
DaTScan2	0.027	2.31	2.08
DaTScan4	0.001	1.4	1.1
<b>MF1-LSTM</b>			
CSF-Total tau	0.007	87.9	46.72
MoCA	$2.16 \times 10^{-17}$	47.5	41.05
SDM	0.005	58.5	41.5
<b>MF2-LSTM</b>			
HVLT-Retention	0.03	0.84	0.83
SDM	0.007	36.61	41.68

- Chi-square test for the categorical features, F-test for the normal continuous features, Kruskal-Wallis test for non-normal continuous features, and Fisher's exact test for sparse features.
- If the  $p$ -value is less than 0.05, a significant group effect is considered for the associated feature.
- Method producing more features with  $p < 0.05$ , is considered as providing a more sensible patient subtyping result.

# Time-Aware LSTM: Parkinson's Disease Patient Subtyping



- Orange color represents the cluster mean which is higher than the total mean of the patients and the shades of blue show lower mean values for the corresponding feature with  $p < 0.05$ .
- PD patients are known to have lower DaTScan SBR values than healthy subjects.

# Drug Similarity Integration Through Attentive Multi-view Graph Auto-Encoders

Tengfei Ma, Cao (Danica) Xiao, Jiayu Zhou, Fei Wang

**IJCAI' 18**

# Drug-Drug Interaction (DDI)

## Combined Use of Drugs



Decrease  
actions of  
drugs

Increase  
actions of  
drugs

Cause  
adverse  
effects



## Facts of DDIs

- Common among patients with complex diseases or comorbidities.
- Hard to observe in clinical testing.
- Affects 15% U.S. population. Cost more than \$177 billion per year in disease management.

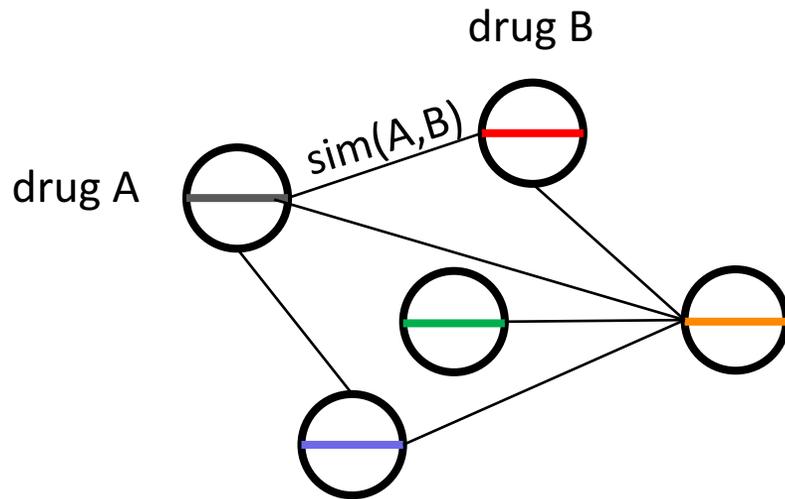
# DDI Prediction based on Multiview Data

- **Drug Features** (database)
  - Label Side Effect (SIDER)
  - Off-Label Side Effect (OFFSIDES)
  - Molecular substructure
  - Drug Indication (MedDRA)
  - .....
- **Assumption:** similar drugs may interact with the same drug.
- **Approach:** We consider each type of feature as a view that has partial correlation with drug similarity, so we aim at integrating similarity metrics across multiple views **for more accurate similarity learning.**

# State-of-the-art and Challenges

- State-of-the-art
  - Nearest neighbor methods (Zhang *et al*, 2015, Zhang *et al*, 2017 )
  - Random walk based methods (Wang *et al*, 2010)
  - Unsupervised iterative methods (Angione *et al*, 2014, Xu *et al*, 2016)
  - Multiple kernel learning (Zhuang *et al*, 2011, McFee *et al*, 2011)
- Challenges
  - the underlying relations of biomedical events are often **nonlinear** and **complex** over all types of features
  - features have **different importance** toward different target outcomes

# Model Overview



- Construct drug similarity graph, where DDIs are node labels
- Graph convolutional network (GCN) based model for node embedding and prediction
- Attention mechanism to integrate similarities from different views

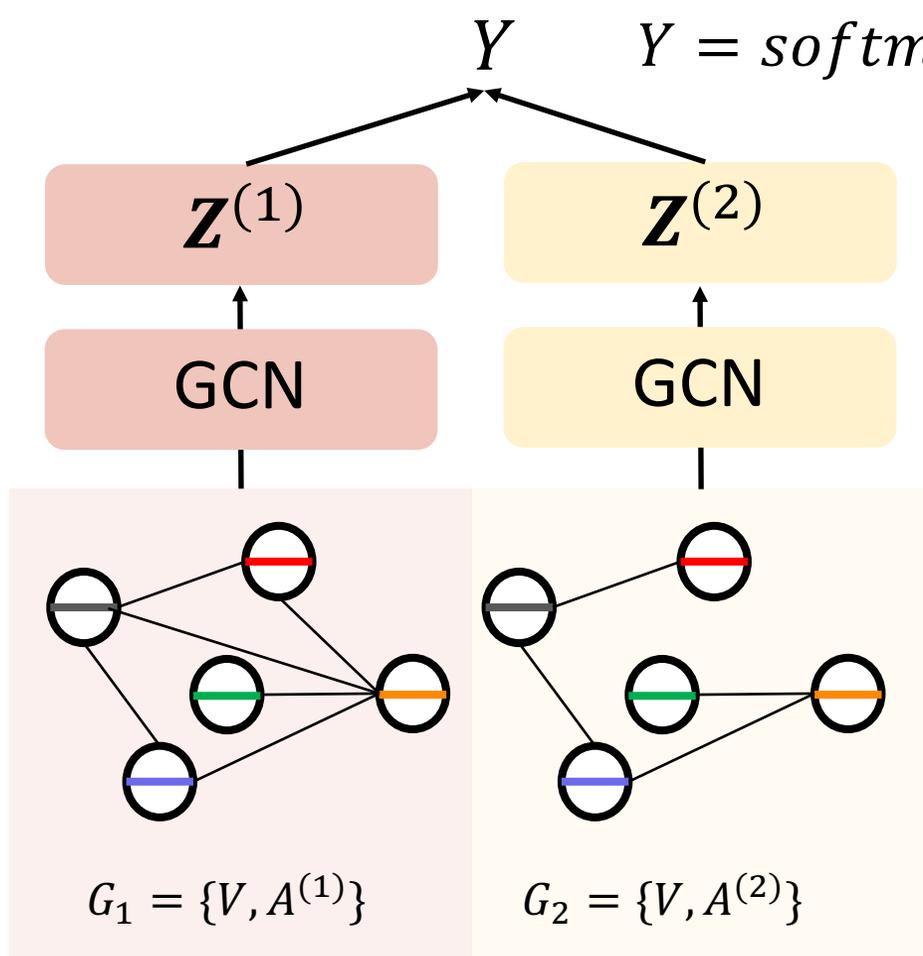
# Recall: GCN (Kipf and Welling, 2016)

- Objective: graph node embedding for arbitrary graphs, distance as (dis-)similarity of local graph structures
- Input
  - Node features  $\mathbf{x}$
  - Adjacency matrix  $\mathbf{A}$  that represents graph structure
- Output: Node embedding  $\mathbf{Z}$
- Method

$$\mathbf{H}^{(l+1)} = f(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)})$$

where  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_N$ ,  $\mathbf{D}$  is a diagonal matrix such that  $D_{ii} = \sum_j \tilde{A}_{ij}$ ,  $\mathbf{W}^{(l)}$  is layer-specific parameter matrix,  $\mathbf{H}^{(l)}$  is node representation of  $l$ th layer.

# Multi-View GCN



## Embedding

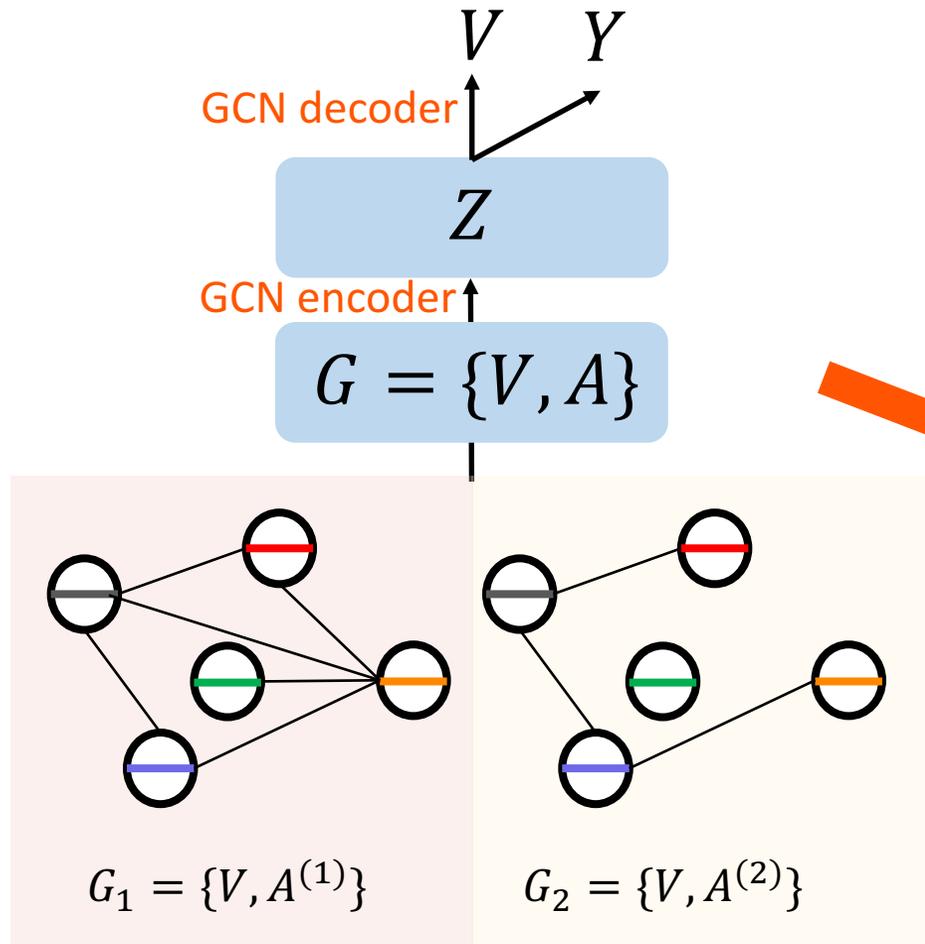
For view  $u \in \{1, \dots, T\}$ , we encode the nodes as

$$\mathbf{Z}^{(u)} = f(\mathbf{X}^{(u)}, \mathbf{A}^{(u)}; \mathbf{W}_1^{(u)})$$

$$= \text{softmax}(\hat{\mathbf{A}}^u \text{ReLU}(\hat{\mathbf{A}}^u \mathbf{X}^{(u)} \mathbf{W}_0^{(u)}) \mathbf{W}_1^{(u)})$$

where  $\hat{\mathbf{A}}^u = \tilde{\mathbf{D}}^{(u)-\frac{1}{2}} \tilde{\mathbf{A}}^{(u)} \tilde{\mathbf{D}}^{(u)-\frac{1}{2}}$ , and  $\mathbf{W}_0^{(u)}$  and  $\mathbf{W}_1^{(u)}$  are weight matrices.

# Attentive Multiview Similarity Fusion with Graph Auto-Encoders (GAE)



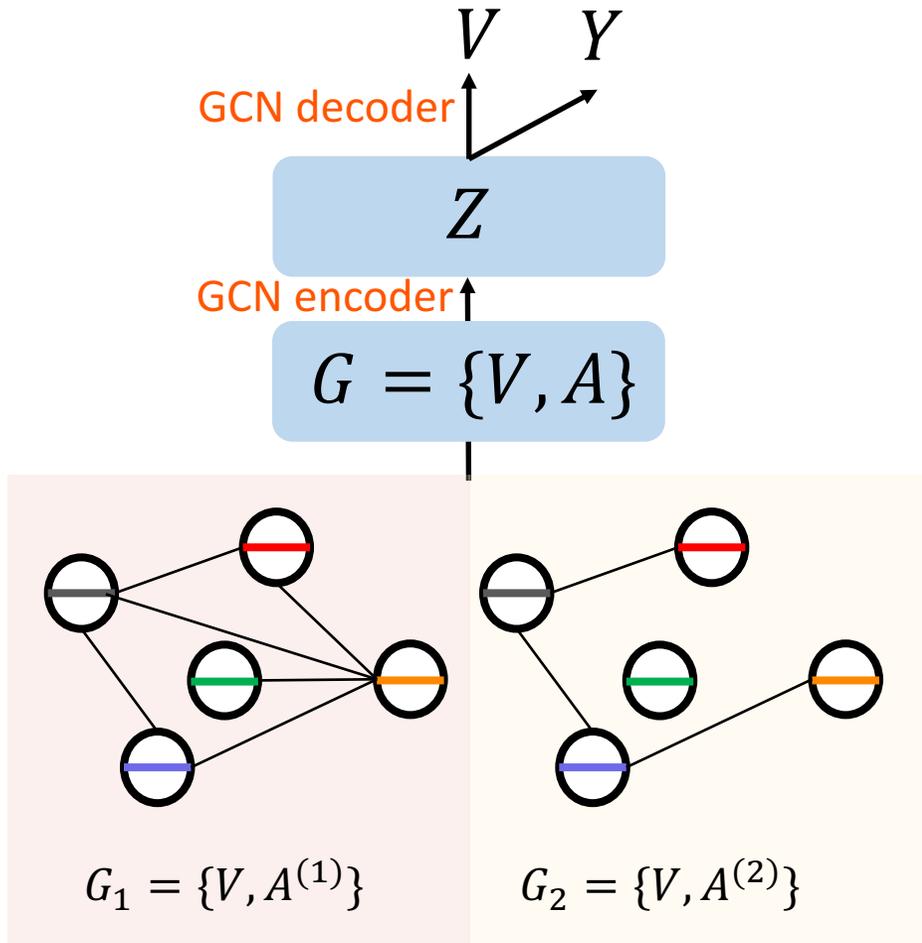
When we do not have labels, we reconstruct  $V$  from  $Z$  to minimize autoencoder loss,  $L_{ed} = \sum |X - X'|^2$ .

$$A = \sum_u \text{diag}(g^u) * A^u$$

$$g'^u = W^u A^u + b^u \xrightarrow{\text{Normalize}} g^u$$

attention weights decided by data and target

# Semi-Supervised Extensions (SemiGAE)



## SemiGAE (for partial labels)

Objective:  $\min L = L_{train} + L_{ed}$

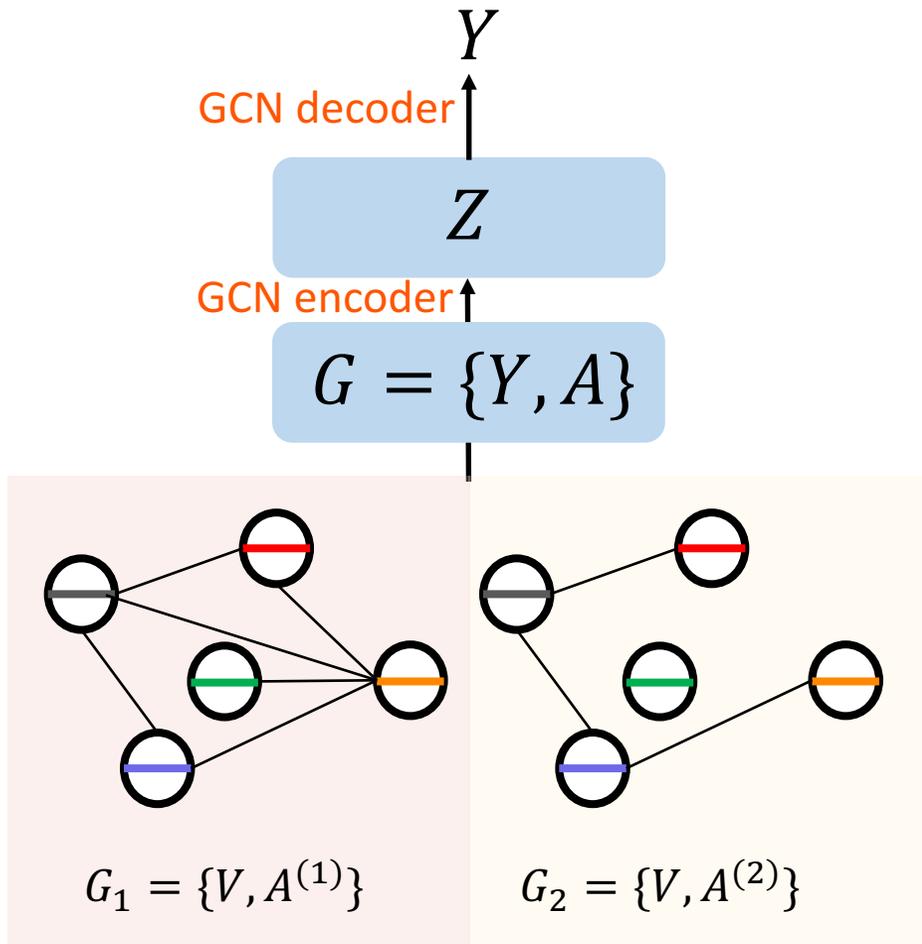
Training loss for labeled data

$$L_{train} = \sum_{y \in Y_{train}} \sum_{y' \in h(Z_{train})} -y \ln y'$$

Auto-encoder loss for all data

$$L_{ed} = \sum |X - X'|^2$$

# Transductive Extensions (TransGAE)



## TransGAE (for lack of node features)

For the case when we do not have node feature

$$\min L = |Y'_{train} - Y_{train}|^2 + |Y'_{test} - Y_{test}|^2 + \mu|Y_{test}|^2$$

where we also treat DDI label as input variable

$$\{Y'_{train}, Y'_{test}\} = f'(f(\{Y_{train}, Y_{test}\}, \hat{A}), \hat{A})$$

# Experiments

Data (Binary)	Dimension
drugs (pairs)	645 (63473)
DDI	1318
label ADR	4192
off-label ADR	10093
Substructure	645 x 1024
Data (Multi)	Dimension
drugs (pairs)	222 (63473)
DDI	1301
indication	1702
CPI	611
TTD	207
Substructure	645 x 582

## Baselines

- Nearest neighbor [Vilar et al. 2012]
- Label Propagation [Zhang et al. 2015]
- Multiple Kernel Learning [Strazar and Curk 2016]
- Basic Multi-view GraphCNN

## Evaluation

- Same strategy as [Zhang et al. 2015]
- Selecting a fixed percentage of drugs randomly and all DDIs associated with these drugs are used for testing
- For the remaining training data, 90%/10% split for training/validation
- Evaluation metrics: ROC-AUC, PR-AUC

# Results

Table 2: Predicting Specific DDI Types (Multiple Outcomes) on Dataset 2.

Using Single View					
	Methods	Test Split (25%)		Test Split (50%)	
		ROC-AUC	PR-AUC	ROC-AUC	PR-AUC
Baselines	NN	$0.627 \pm 0.043$	$0.594 \pm 0.078$	$0.594 \pm 0.033$	$0.554 \pm 0.061$
	LP	$0.773 \pm 0.025$	<b><math>0.670 \pm 0.052</math></b>	$0.747 \pm 0.028$	<b><math>0.650 \pm 0.053</math></b>
	GraphCNN	$0.738 \pm 0.047$	$0.594 \pm 0.080$	$0.698 \pm 0.090$	$0.583 \pm 0.102$
Proposed	SemiGAE	<b><math>0.798 \pm 0.029</math></b>	$0.661 \pm 0.059$	<b><math>0.784 \pm 0.028</math></b>	$0.649 \pm 0.059$
	TransGAE	$0.790 \pm 0.028$	$0.661 \pm 0.068$	$0.770 \pm 0.031$	$0.633 \pm 0.080$
Using Multiple Views					
Baselines	LP	$0.774 \pm 0.025$	$0.672 \pm 0.052$	$0.748 \pm 0.028$	$0.653 \pm 0.055$
	GraphCNN	$0.601 \pm 0.067$	$0.526 \pm 0.120$	$0.578 \pm 0.067$	$0.526 \pm 0.108$
	MKL	$0.766 \pm 0.030$	$0.650 \pm 0.061$	$0.724 \pm 0.026$	$0.586 \pm 0.066$
Proposed	AttSemiGAE	<b><math>0.802 \pm 0.029</math></b>	<b><math>0.678 \pm 0.060</math></b>	<b><math>0.786 \pm 0.030</math></b>	<b><math>0.662 \pm 0.064</math></b>
	AttTransGAE	$0.782 \pm 0.026$	$0.670 \pm 0.058$	$0.764 \pm 0.025$	$0.652 \pm 0.061$

# Analysis of Interpretability

DDI Type	AUC	Attention Weights			
		Chem.	indi.	TTDS	CPI
Chest Pain	0.772	0.151	0.303	0.144	0.402
Insomnia	0.755	0.380	0.261	0.078	0.291
indication	0.774	0.117	0.301	0.283	0.299

## Results

Views “indication” and “CPI” receive high weights for the ADR “Chest pain”.



## Clinical Evidence

Many DDI cases of chest pain are due to particular drug overdose [Nachimuthu et al., 2012]. For example, the co-use of many medications that treat depression can cause overdose and prolong the QT interval via CPI, and eventually cause chest pain.

# Data Augmentation

- MedGAN (generate EHR statistics)
- Real-valued (medical) Time Series Generation with Recurrent Conditional GANs (NIPS ML4H 17)
- Recent Trend in Molecular Graph Generation for Drug Discovery

# Generating Multi-label Discrete Electronic Health Records using Generative Adversarial Networks

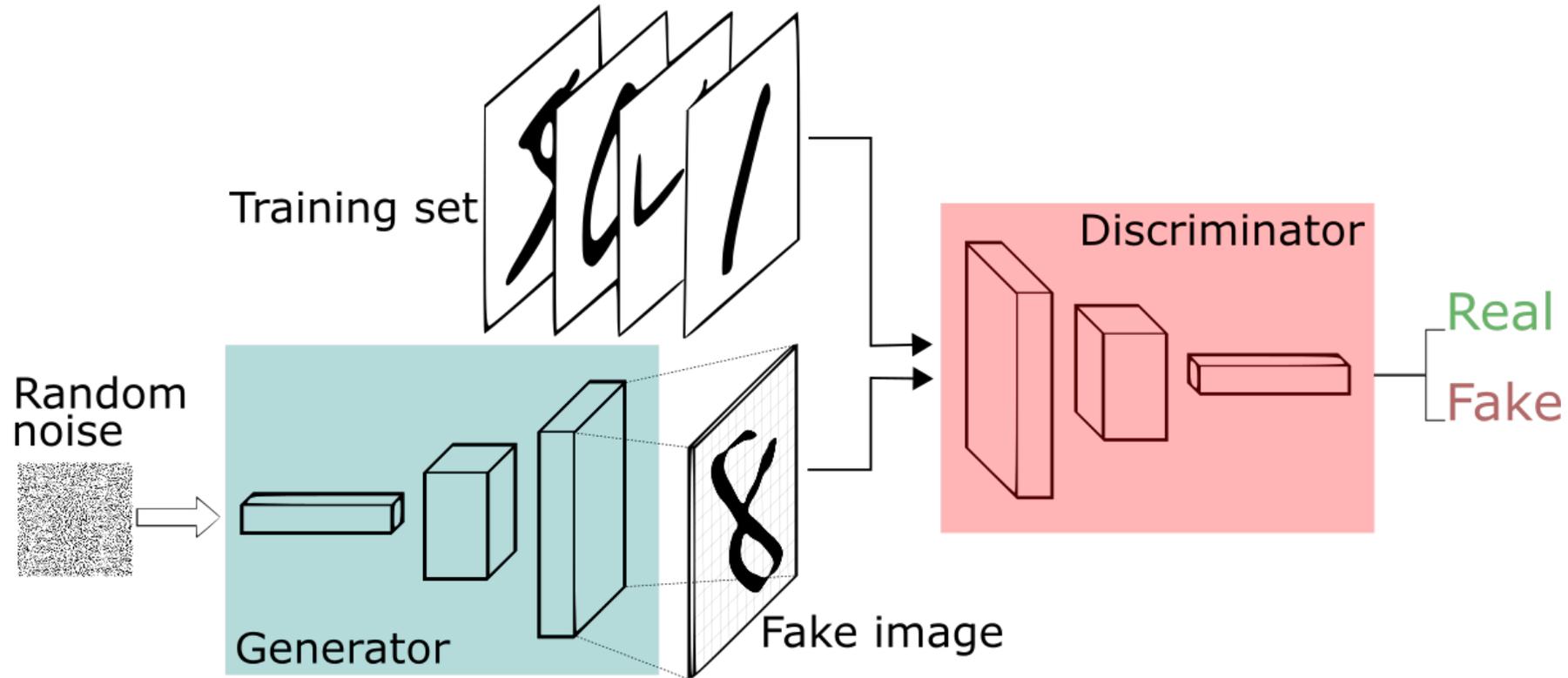
Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, Jimeng Sun

**MLHC 2017**

# medGAN: Background

- Generating Multi-label Discrete Electronic Health Records using Generative Adversarial Networks
  - Choi, Edward, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun, MLHC 2017
- Generate synthetic patient records that are
  - Similar (statistically) to the real records
  - Does not divulge individual patient information

# medGAN: Generative Adversarial Networks



<https://sthalles.github.io/intro-to-gans/>

# medGAN: Patient Record

- Real patient dataset
  - Patients' records are aggregated over time
  - Discrete values (count, binary)

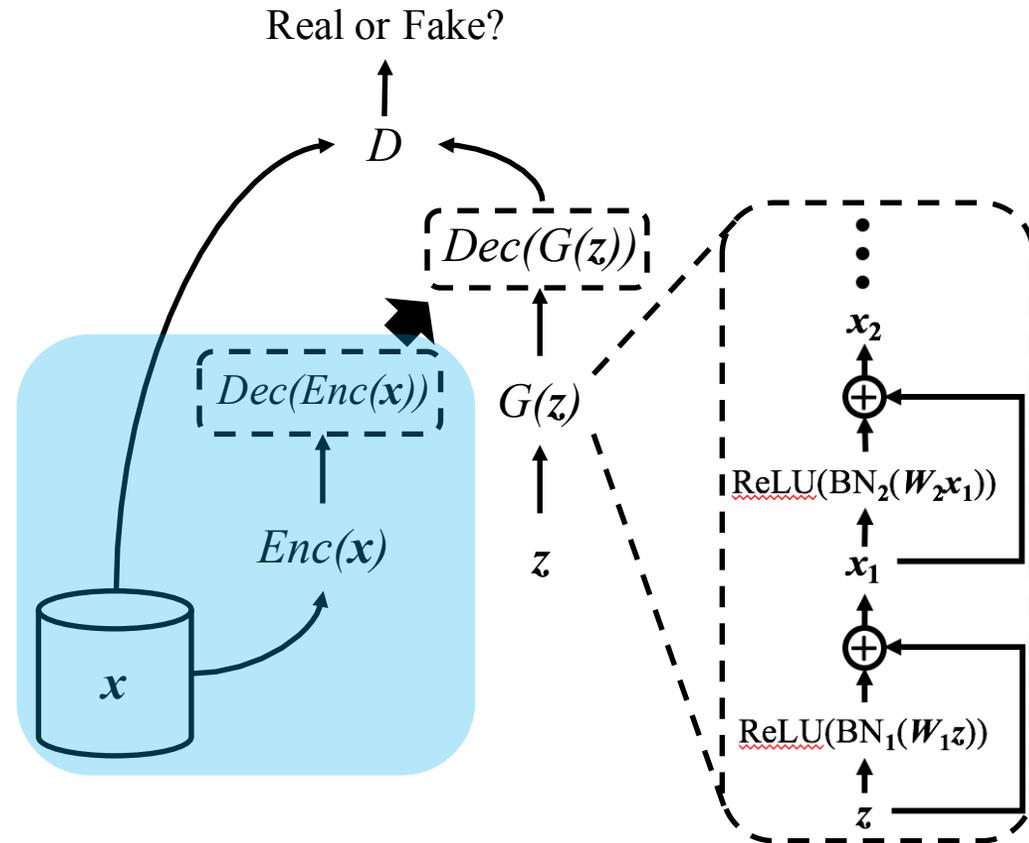
	Hypertension	Heart failure	Diabetes	Kidney failure				
Patient 1	3	0	2	1	0	0	0	...
Patient 2	0	1	0	5	0	4	0	...
Patient 3	1	0	2	2	0	0	6	...

# medGAN Architecture

Pre-train Autoencoder (ex: binary variable)

$$\frac{1}{m} \sum_{i=0}^m \mathbf{x}_i \log \mathbf{x}'_i + (1 - \mathbf{x}_i) \log(1 - \mathbf{x}'_i)$$

where  $\mathbf{x}'_i = Dec(Enc(\mathbf{x}_i))$

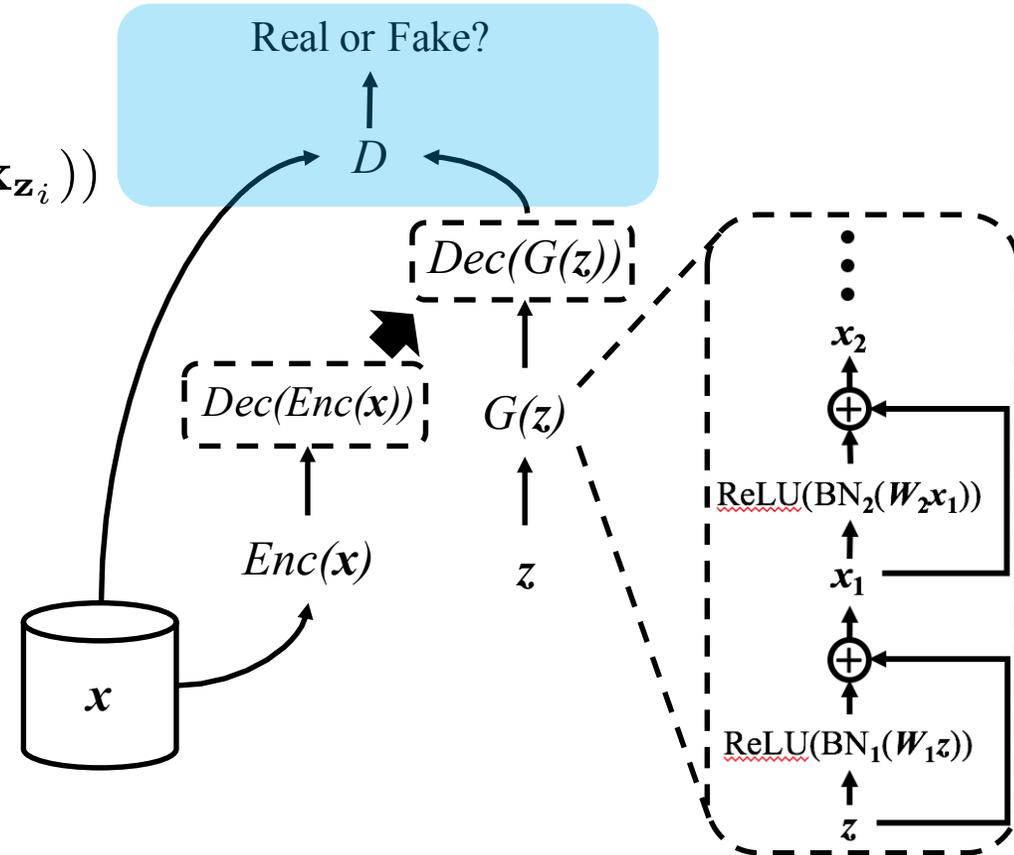


# medGAN Architecture

Update Discriminator parameters

$$\theta_d \leftarrow \theta_d + \alpha \nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \log D(\mathbf{x}_i) + \log(1 - D(\mathbf{x}_{z_i}))$$

where  $\mathbf{x}_{z_i} = Dec(G(\mathbf{z}_i))$

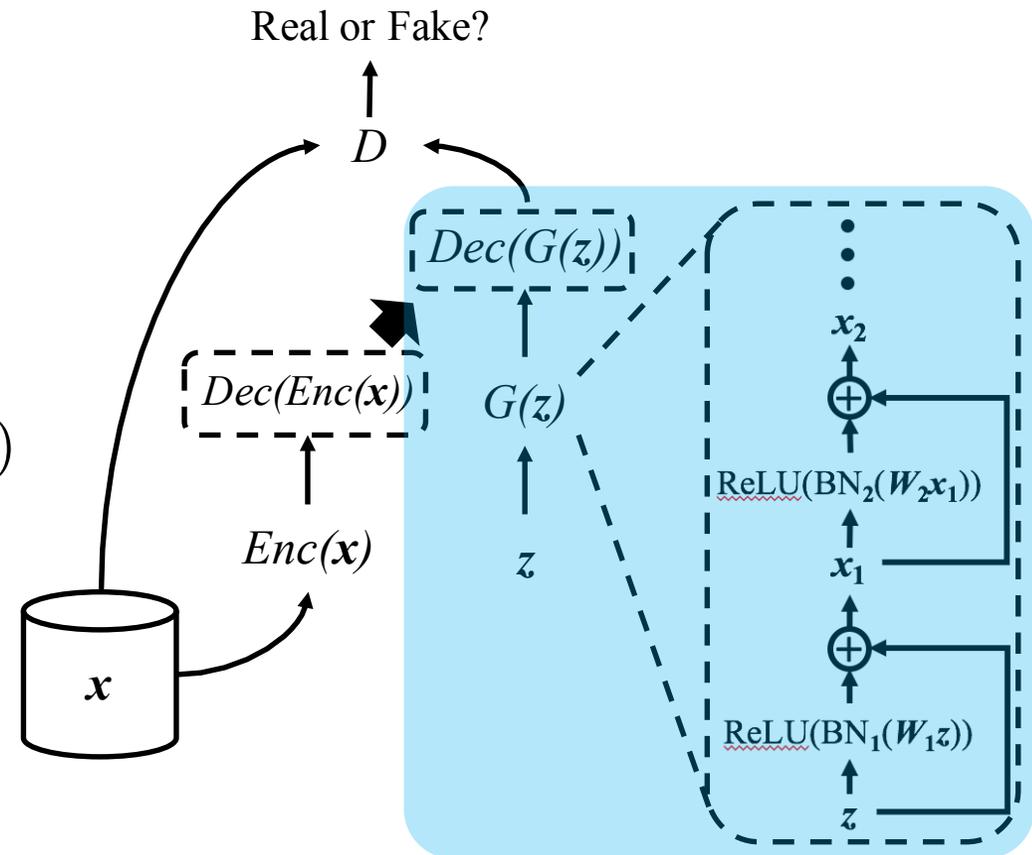


# medGAN Architecture

Update Generator parameters

$$\theta_{g,dec} \leftarrow \theta_{g,dec} + \alpha \nabla_{\theta_{g,dec}} \frac{1}{m} \sum_{i=1}^m \log D(\mathbf{x}_{z_i})$$

where  $\mathbf{x}_{z_i} = Dec(G(\mathbf{z}_i))$

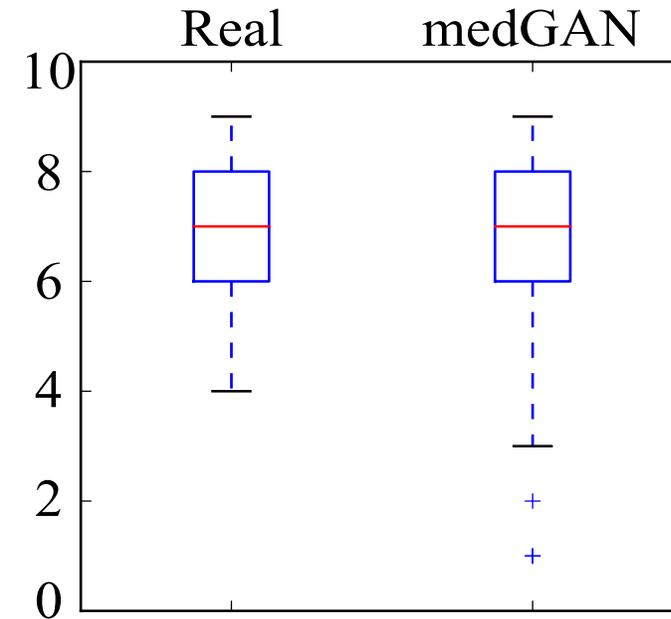
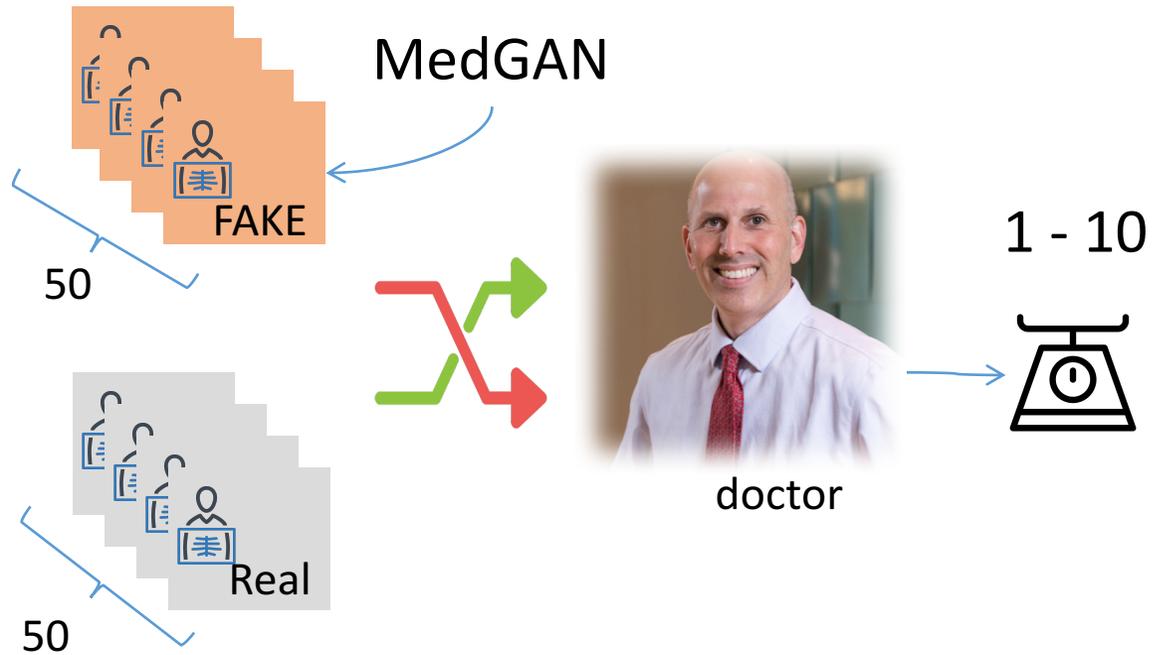


# medGAN: Data

- Data
  - 260K patients from Sutter Health
  - Patient records over 10 years
  - 615 variables
    - Diagnosis/medication/procedure codes
- Binarized the count values
  - Either the code occurred, or did not occur

# medGAN: Result

- Qualitative evaluation by a medical expert



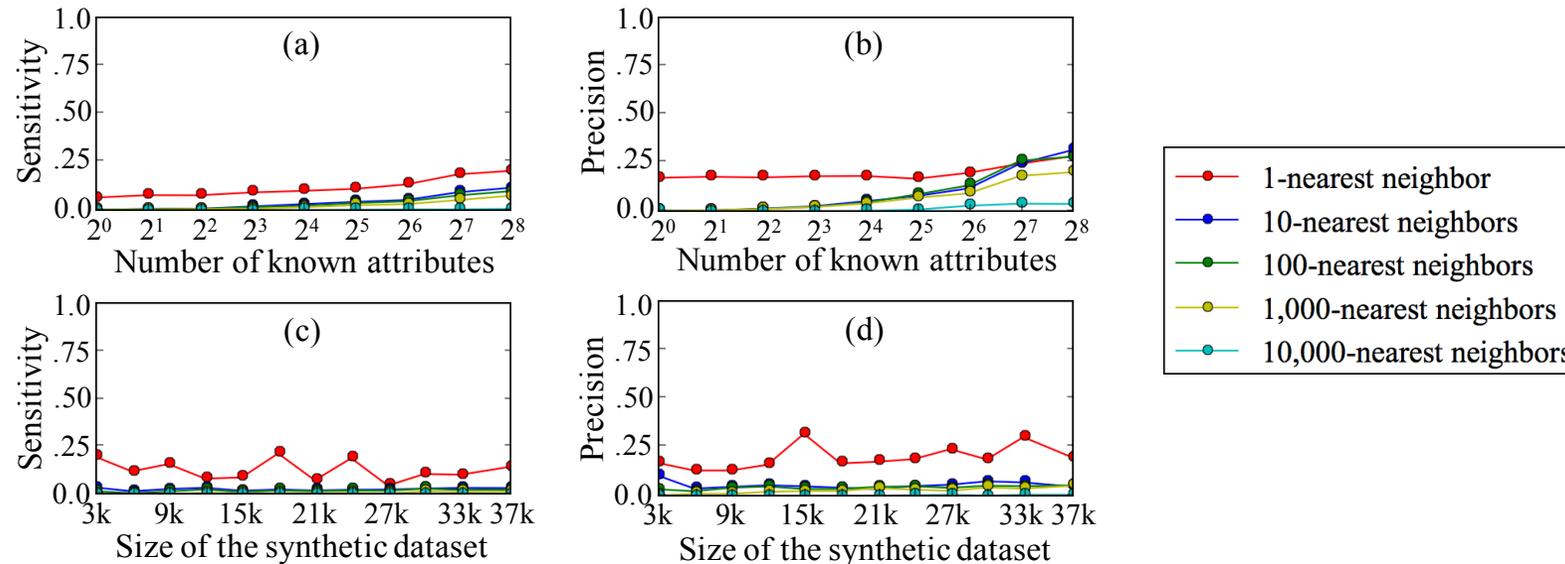
Indistinguishable to a human doctor except a few outliers!

# medGAN: Privacy

- Attribute disclosure attack
  - When an attacker knows subset of patients' info, can they know more by looking at synthetic records?

# medGAN: Privacy

- Attribute disclosure attack



- The attacker who knows 1% of the target patient's attributes (8-16 attributes)
  - Can correctly estimate at most 10% of positive unknown attributes
- Size of the synthetic data has little influence on the effectiveness of the attack

Attacker cannot gain useful knowledge on unknown attributes!

# Real-valued (medical) Time Series Generation with Recurrent Conditional GANs

Esteban, Cristóbal, Stephanie L. Hyland, and Gunnar Rätsch

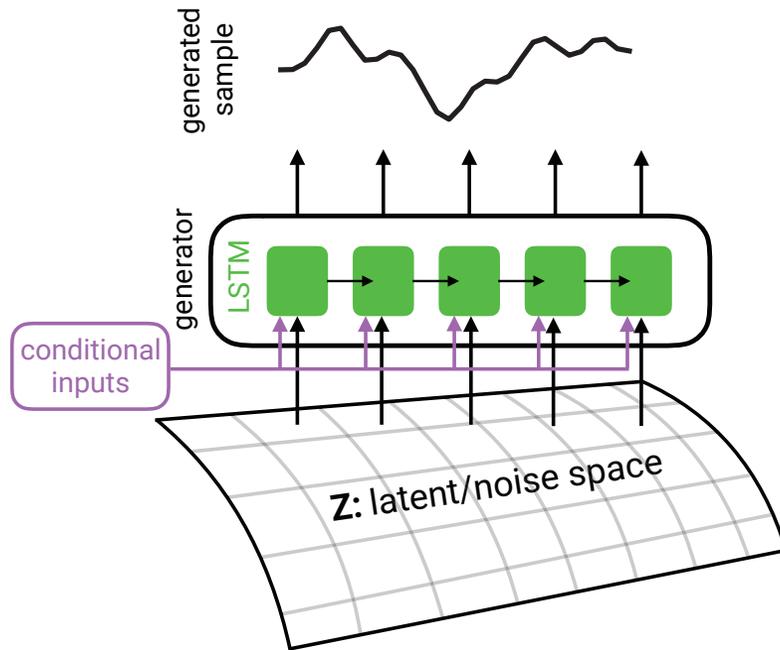
**NIPS ML4H 17**

# RCGAN

- Real-valued (medical) time series generation with recurrent conditional GANs.
  - Esteban, Cristóbal, Stephanie L. Hyland, and Gunnar Rätsch, NIPS ML4H Workshop 2017
- Generate a continuous-valued multi-variate time-series
  - Values such as measurements from ICU

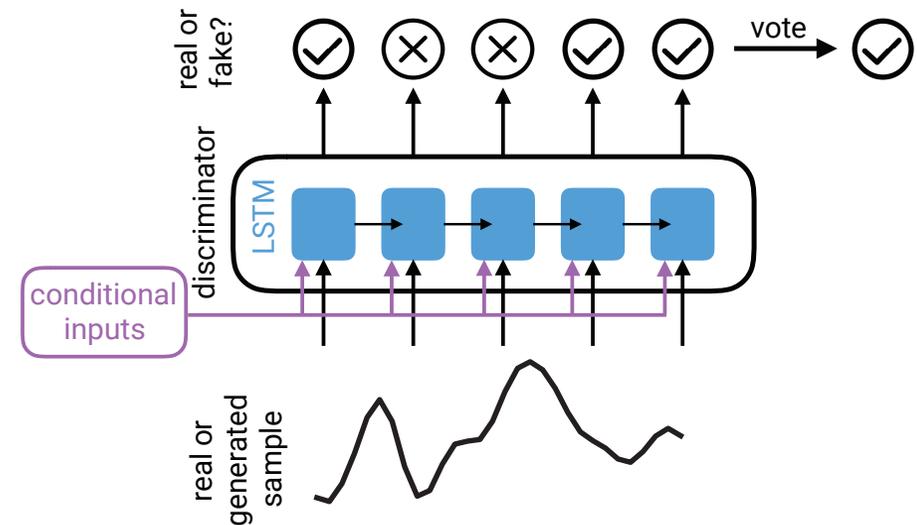
# RCGAN: Model

- Generating time-series using GAN



Generator

$$G_{\text{loss}}(Z_n) = -\text{CE}(\text{RNN}_D(\text{RNN}_G(Z_n)), \mathbf{1})$$



Discriminator

$$D_{\text{loss}}(X_n, \mathbf{y}_n) = -\text{CE}(\text{RNN}_D(X_n), \mathbf{y}_n)$$

# RCGAN: Data

- Dataset
  - 17,693 patients from Philips eICU dataset
- Features
  - Four most frequently recorded variables by bed-side monitors
  - Oxygen saturation (SpO<sub>2</sub>), heart rate (HR), respiratory rate (RR), mean arterial pressure (MAP)
  - 16 measurement per each variable (16-step time-series)
    - Measurements for a 4-hour period
  - Exclude patients with missing data

# RCGAN: Result

- Evaluation of the synthetic data
  - Train Random Forest on synthetic training data
  - Test the trained RF on real test data
- Prediction tasks
  - Based on 4-hours observation (16-timesteps), predict future binary label
  - Will {SpO2, HR, RR, MAP} be high/low in the next hour?

		<i>SpO2 &lt; 95</i>	HR < 70	<i>HR &gt; 100</i>
AUROC	real	$0.9587 \pm 0.0004$	$0.9908 \pm 0.0005$	$0.9919 \pm 0.0002$
	TSTR	$0.88 \pm 0.01$	$0.96 \pm 0.01$	$0.95 \pm 0.01$
AUPRC	real	$0.9059 \pm 0.0005$	$0.9855 \pm 0.0002$	$0.9778 \pm 0.0002$
	TSTR	$0.66 \pm 0.02$	$0.90 \pm 0.02$	$0.84 \pm 0.03$
	random	0.16	0.26	0.18

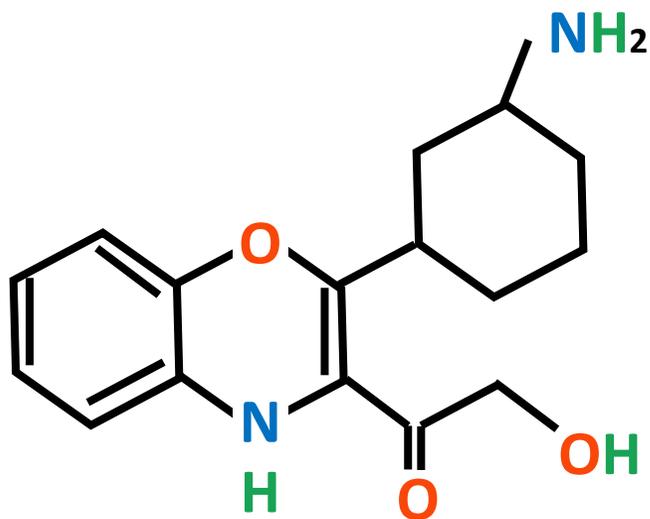
		<i>RR &lt; 13</i>	RR > 20	MAP < 70	MAP > 110
AUROC	real	$0.9735 \pm 0.0001$	$0.963 \pm 0.001$	$0.9717 \pm 0.0001$	$0.960 \pm 0.001$
	TSTR	$0.86 \pm 0.01$	$0.84 \pm 0.02$	$0.875 \pm 0.007$	$0.87 \pm 0.04$
AUPRC	real	$0.9557 \pm 0.0002$	$0.891 \pm 0.001$	$0.9653 \pm 0.0001$	$0.8629 \pm 0.0007$
	TSTR	$0.73 \pm 0.02$	$0.50 \pm 0.06$	$0.82 \pm 0.02$	$0.42 \pm 0.07$
	random	0.26	0.1	0.39	0.05

# RCGAN: Conclusion

- Is RCGAN just memorizing?
  - Are there privacy risks?
- Use Maximum Mean Discrepancy (MMD) for evaluation
  - MMD: measure of the distance between two distributions
- Null hypothesis: RCGAN has **not** memorized the training data
  - $\text{MMD}(\text{synthetic data, test data}) < \text{MMD}(\text{synthetic data, train data})$
- P-value: 0.40+-0.05
  - Cannot reject the null hypothesis
- Synthetic samples do not look more similar to the training set than they do to the test set

# Recent Trend in Variational Autoencoders based Molecular Graph Generation for Drug Discovery

# Molecular Graph Generation: Molecule Encoding



The simplified molecular-input line-entry system (SMILES) is a format of encoding molecules in text strings.

**OCC(=O)C=1Nc3ccccc3OC=1C2CCC(N)C2**

Consider a molecule as a sequential data and apply NLP techniques.

# Molecular Graph Generation: Benchmark Data

	# Graphs	# Nodes	# Node Types	# Edge Types
QM9 (molecule)	134k	9	4	3
ZINC (molecule)	250k	38	9	3

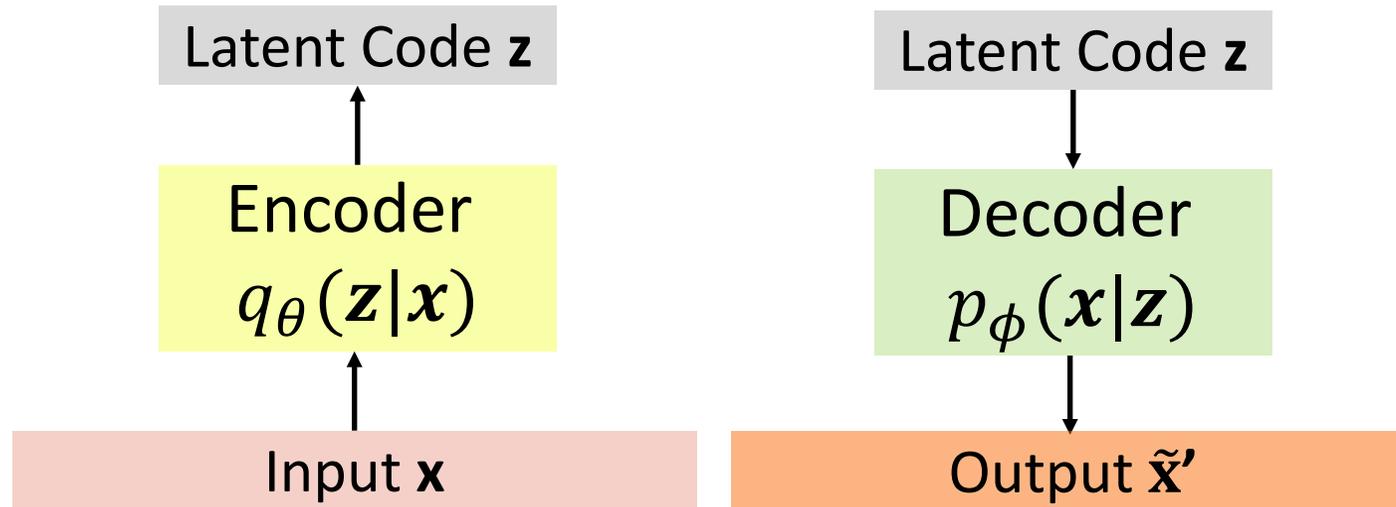
**QM9:** L. Ruddigkeit, R. van Deursen, L. C. Blum, J.-L. Reymond, Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17, *J. Chem. Inf. Model.* 52, 2864–2875, 2012.

**ZINC:** Irwin, Sterling, Mysinger, Bolstad and Coleman, *J. Chem. Inf. Model.* 2012 DOI: 10.1021/ci3001277.

# Molecular Graph Generation: Challenges & Opportunities

- **Goal:** Train generative models to construct more complex, discrete data types.
- **Challenge:** 1. Formal Languages is very strict 2. Small changes in output leads to syntax error
- **Opportunities:** 1. Syntax is context free 2. Grammar is known and fixed 3. Parses are unique

# Tools: Variational AutoEncoder (Kingma et al., 13; Rezende et al., 14)

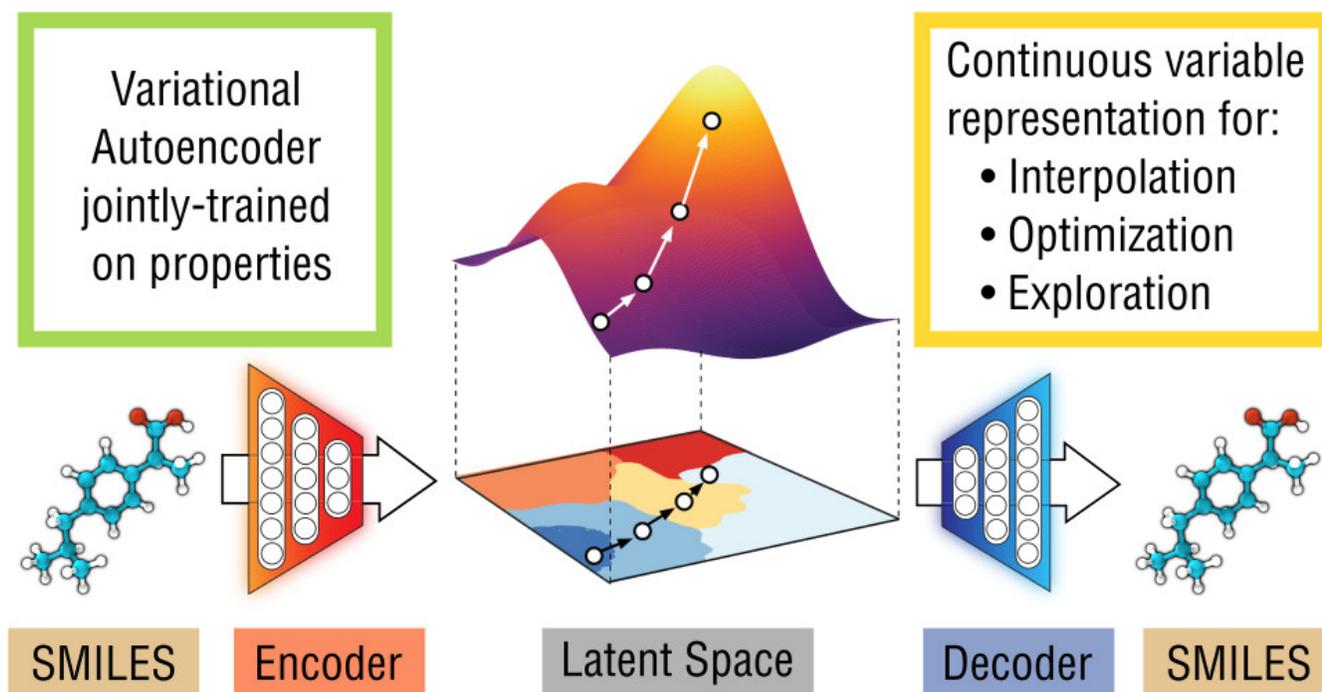


$$L(\theta, \phi) = \underbrace{-E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log(p_{\phi}(\mathbf{x}|\mathbf{z}))]}_{\text{reconstruction loss}} + \underbrace{KL(q_{\theta}(\mathbf{z}|\mathbf{x}), p_{\phi}(\mathbf{z}))}_{\text{penalty for information loss}}$$

reconstruction loss

penalty for  
information loss

# Molecular Graph Generation: CVAE (Rafael Gómez-Bombarelli et al., 2016)

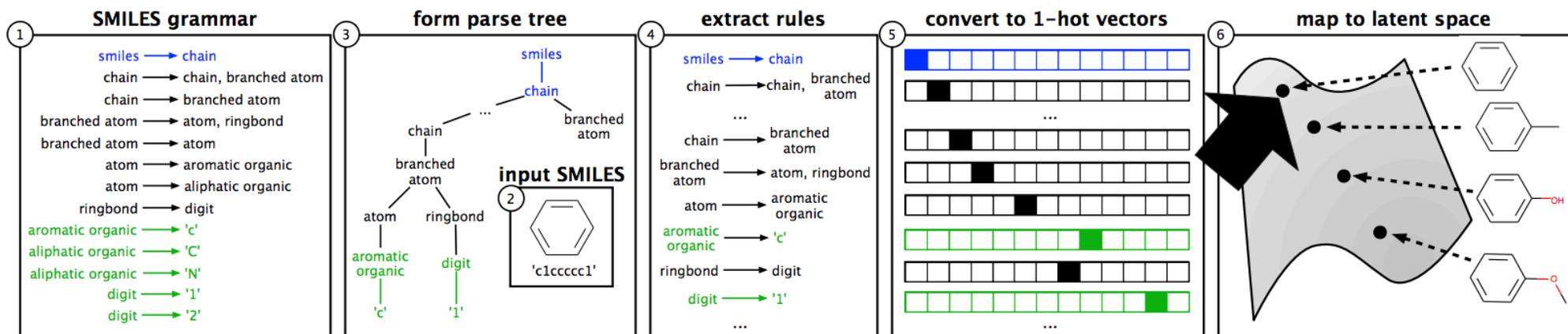


- Encode and decode molecules represented as SMILES with VAE.
- Latent representations could be used to learn in semi-supervised setting.
- The learned model can be used to find molecules with desired property by optimizing representation in latent space and decoding it.

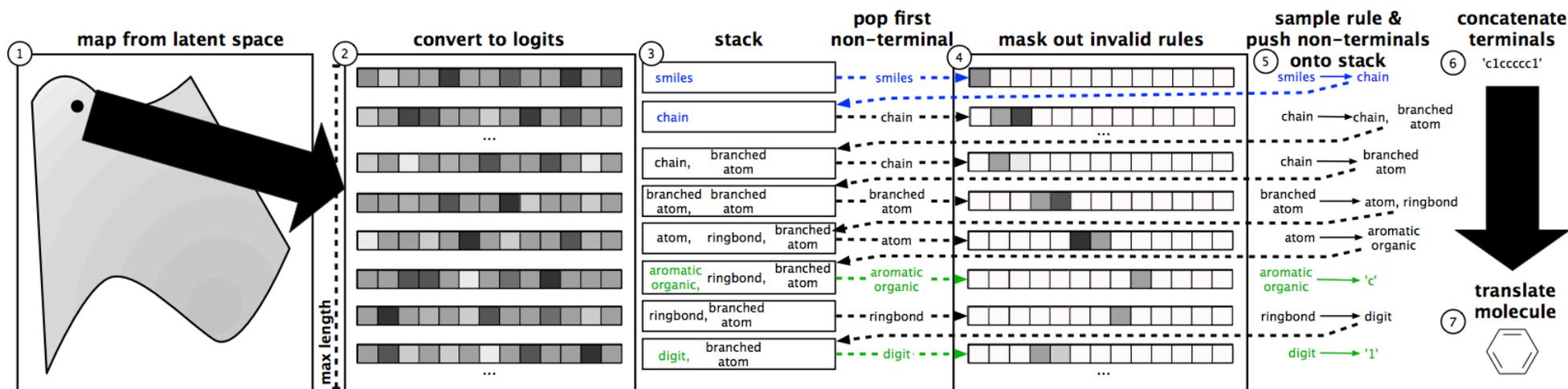
Rafael Gómez-Bombarelli *et al.*, Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules., ACS Cent. Sci., 2018, 4 (2), pp 268–276

# Molecular Graph Generation: Grammar VAE (GVAE, Kusnar et al., 2017)

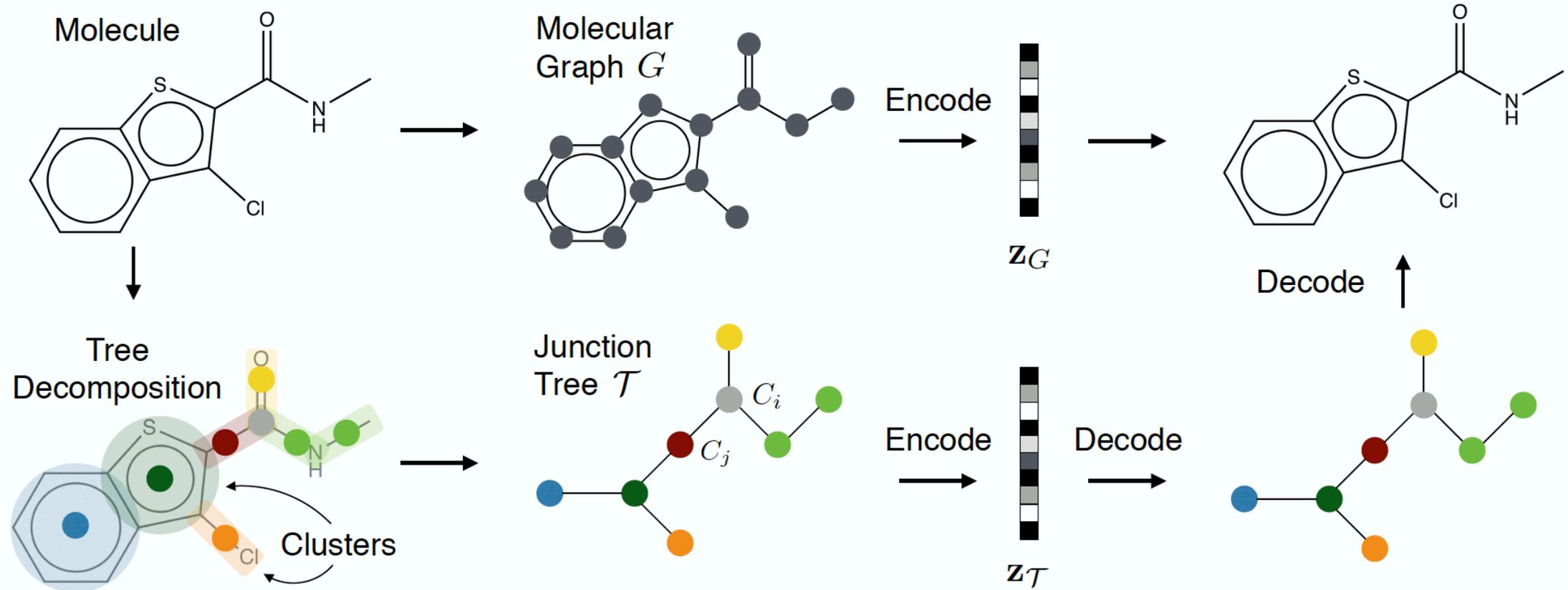
Encode



Decode



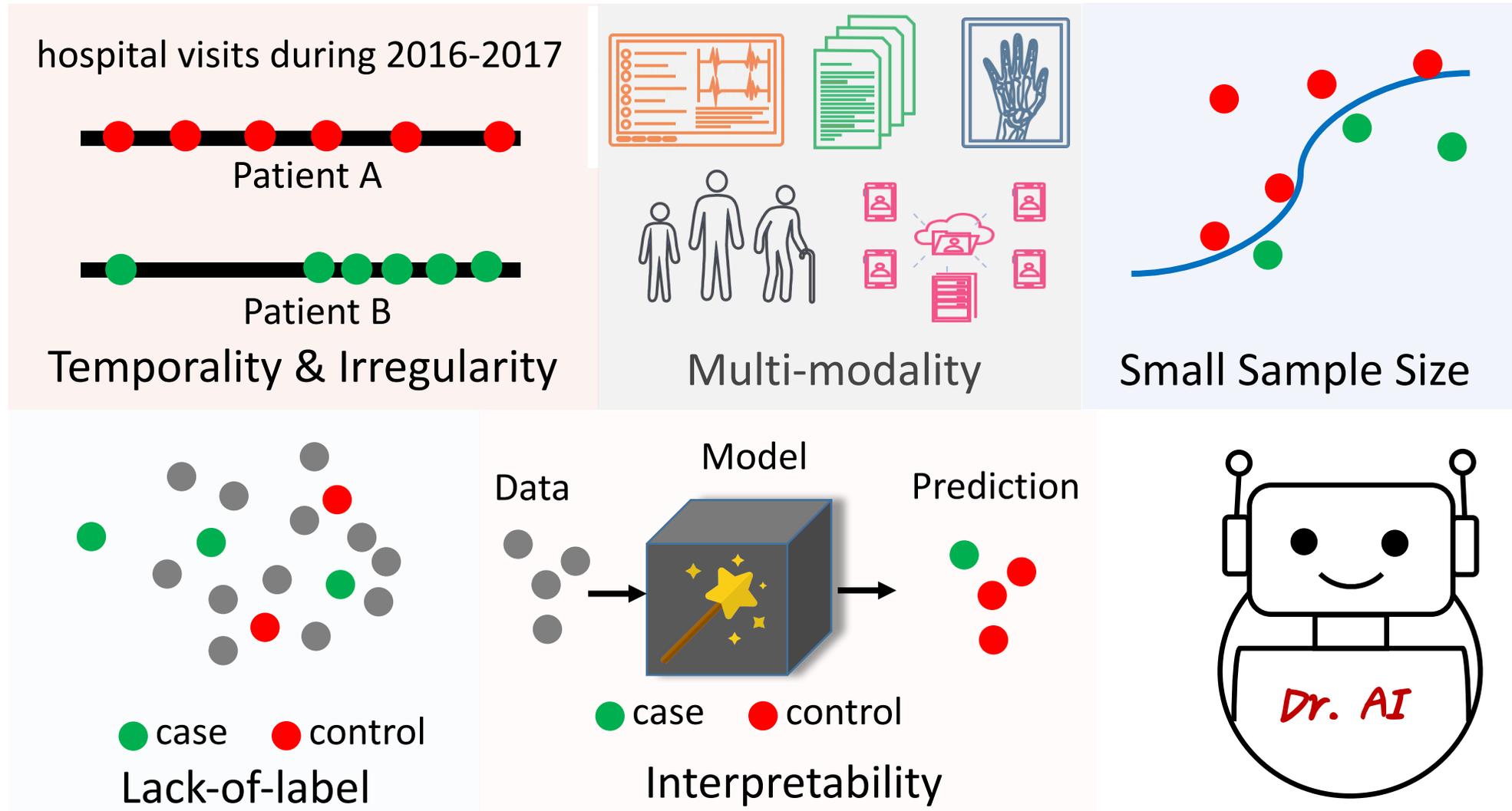
# Molecular Graph Generation: Junction Tree VAE (Jin et al., 2018)



# Agenda

- Background (30 min)
  - healthcare data
  - analytical tasks
  - why deep learning models?
  - deep learning architectures
- Success of Deep Learning in Computational Healthcare (2 hr)
  - Medical Classification
  - Sequential Prediction
  - Concept Embedding
  - Data Augmentation
- Open Challenges
- Q&A

# Challenges of Deep Learning for Healthcare Applications



# What's next?

Modeling heterogeneous data sources



Clinical notes



-Omic data



sensor



Medical imaging

Model interpretation



More complex output



Clinical  
question & answer

# Deep Learning for Healthcare

*Jimeng Sun, Cao (Danica) Xiao, Edward Choi*

*August 7, 2018*