

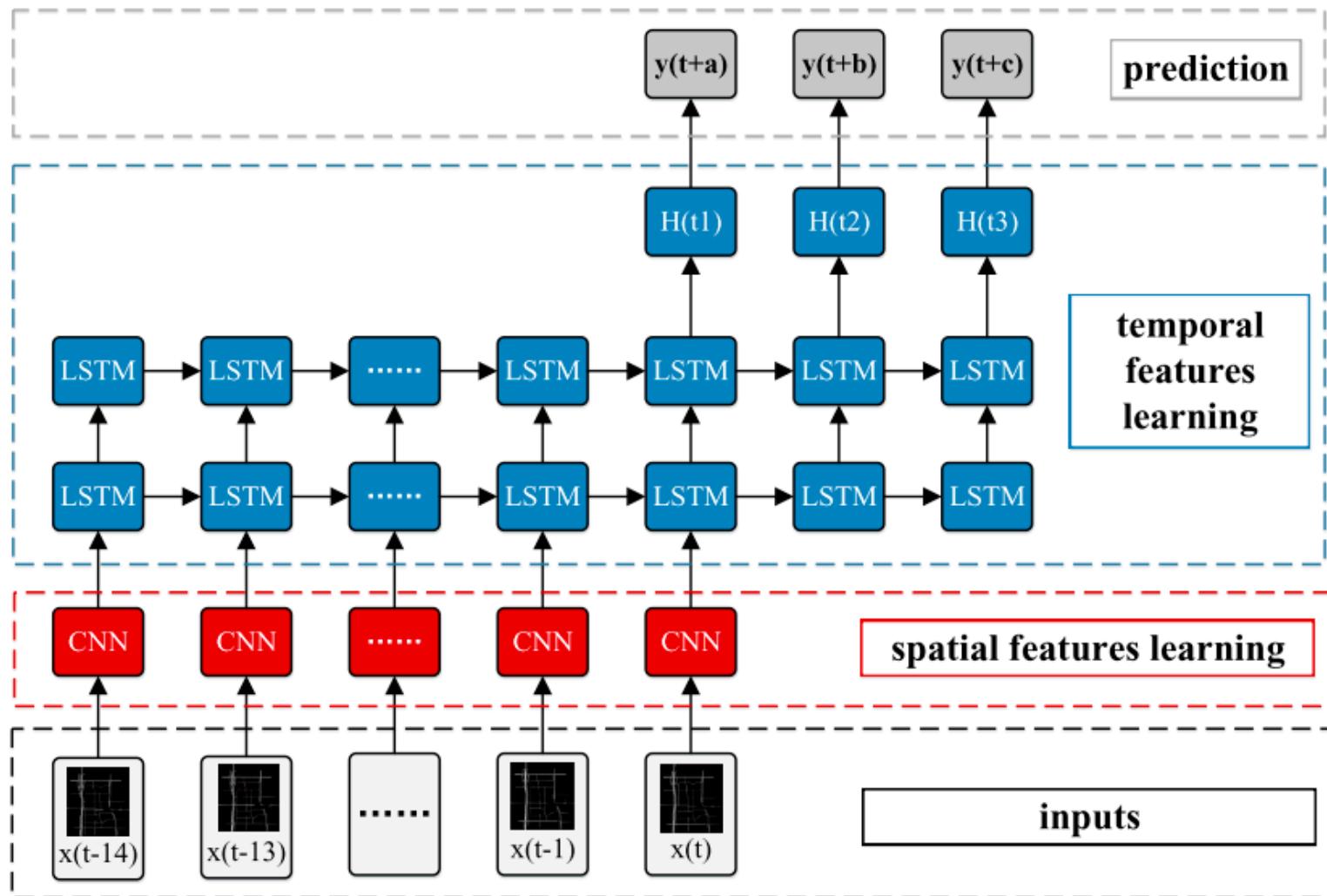
Lecture 4: Embeddings, ConvLSTM, Attention, Graph Neural Networks

Peyman Noursalehi

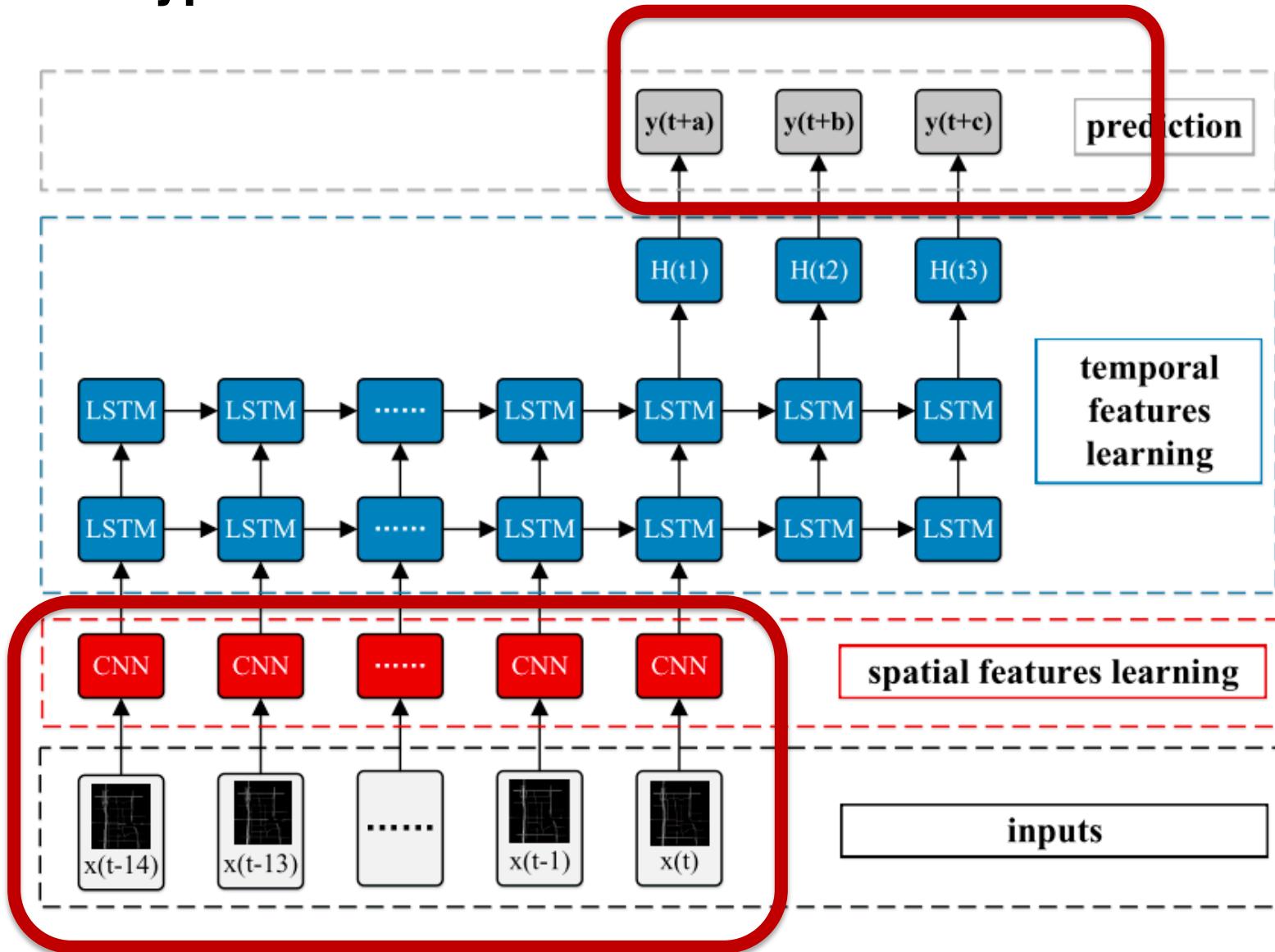
Summary so far

- CNN is the perfect tool for capturing spatial dependencies
- RNN is the perfect tool for capturing temporal dependencies
 - LSTM
 - GRU
- Use CNN in the initial layers to capture spatial dependencies
 - Feature extraction
- Use RNN on features extracted by CNN

A new type of architecture

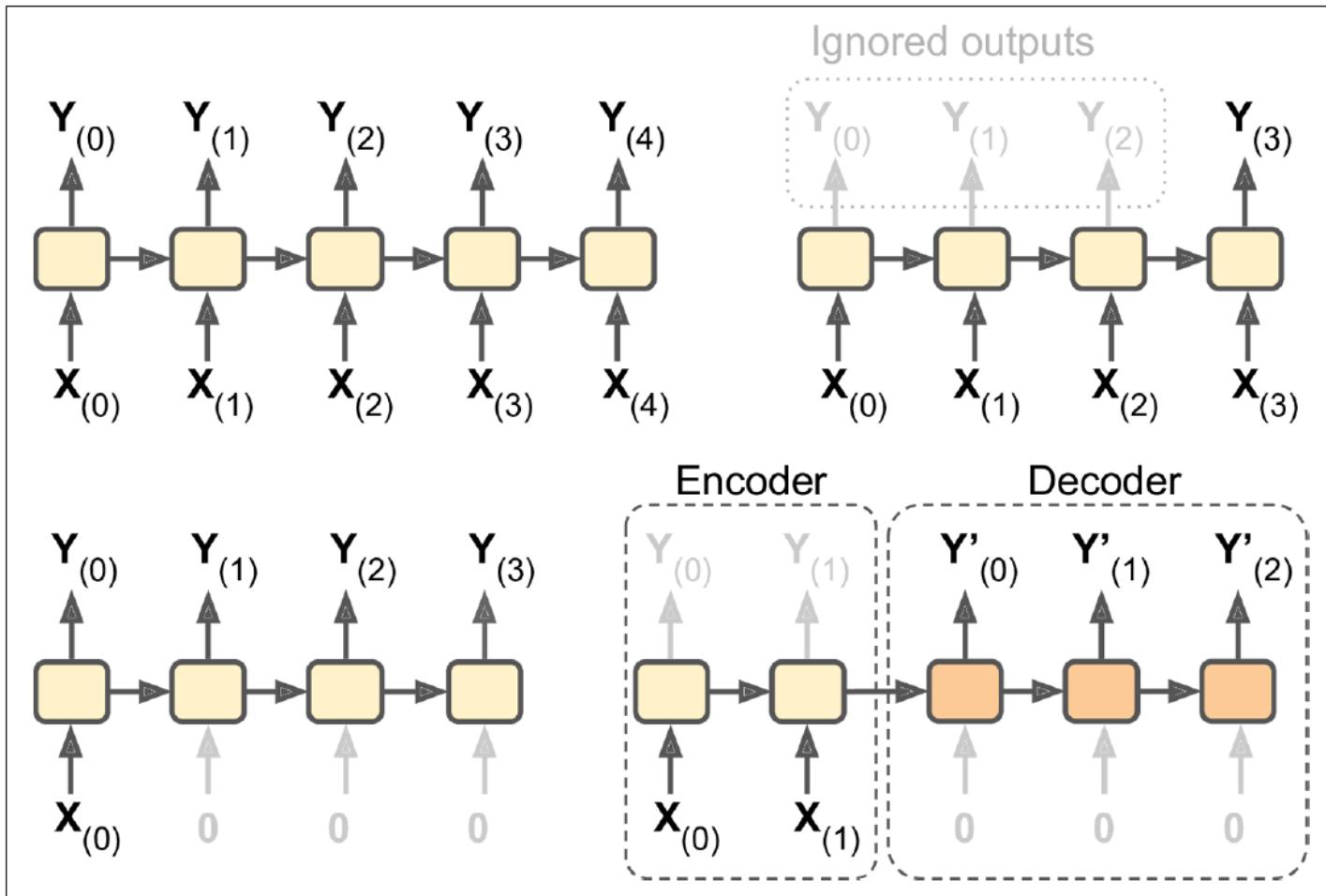


A new type of architecture



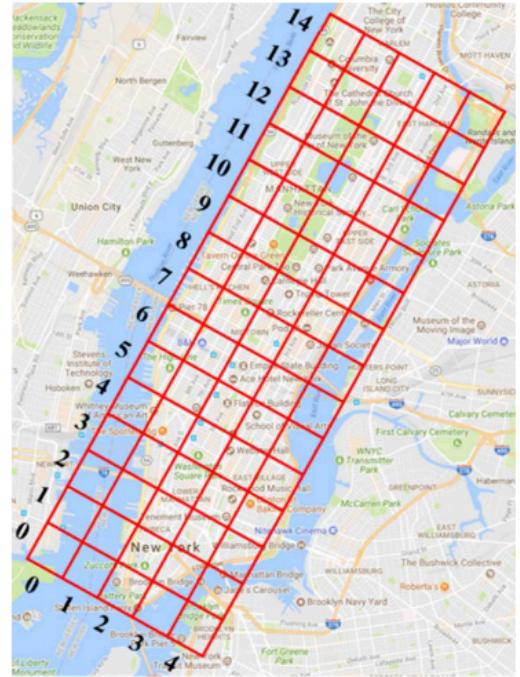
Encoder-decoder architectures

- Many possible formulations
- Encoder-decoder



Taxi demand prediction

- City is divided into N non-overlapping zones
- Predict the number of taxi requests per zone



Deep Multi-View Spatial-Temporal Network for Taxi Demand Prediction

Huaxiu Yao*, Fei Wu

Pennsylvania State University
{huaxiuyao, fxw133}@ist.psu.edu

Jintao Ke*

Hong Kong University of Science
and Technology
jke@connect.ust.hk

Xianfeng Tang

Pennsylvania State University
xianfeng@ist.psu.edu

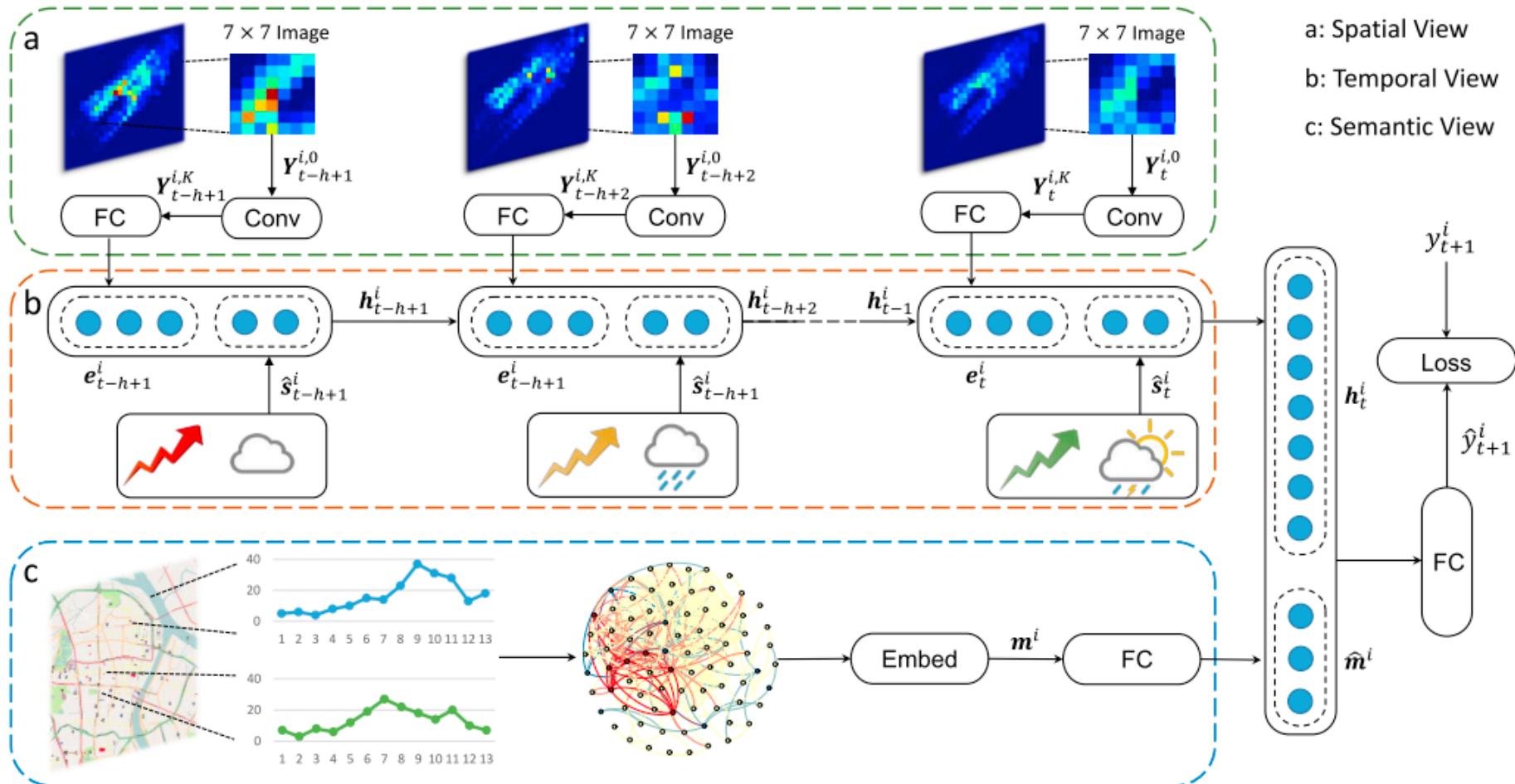
Yitian Jia, Siyu Lu, Pinghua Gong, Jieping Ye

Didi Chuxing
{jiayitian, lusiyu, gongpinghua, yejieping}@didichuxing.com

Zhenhui Li

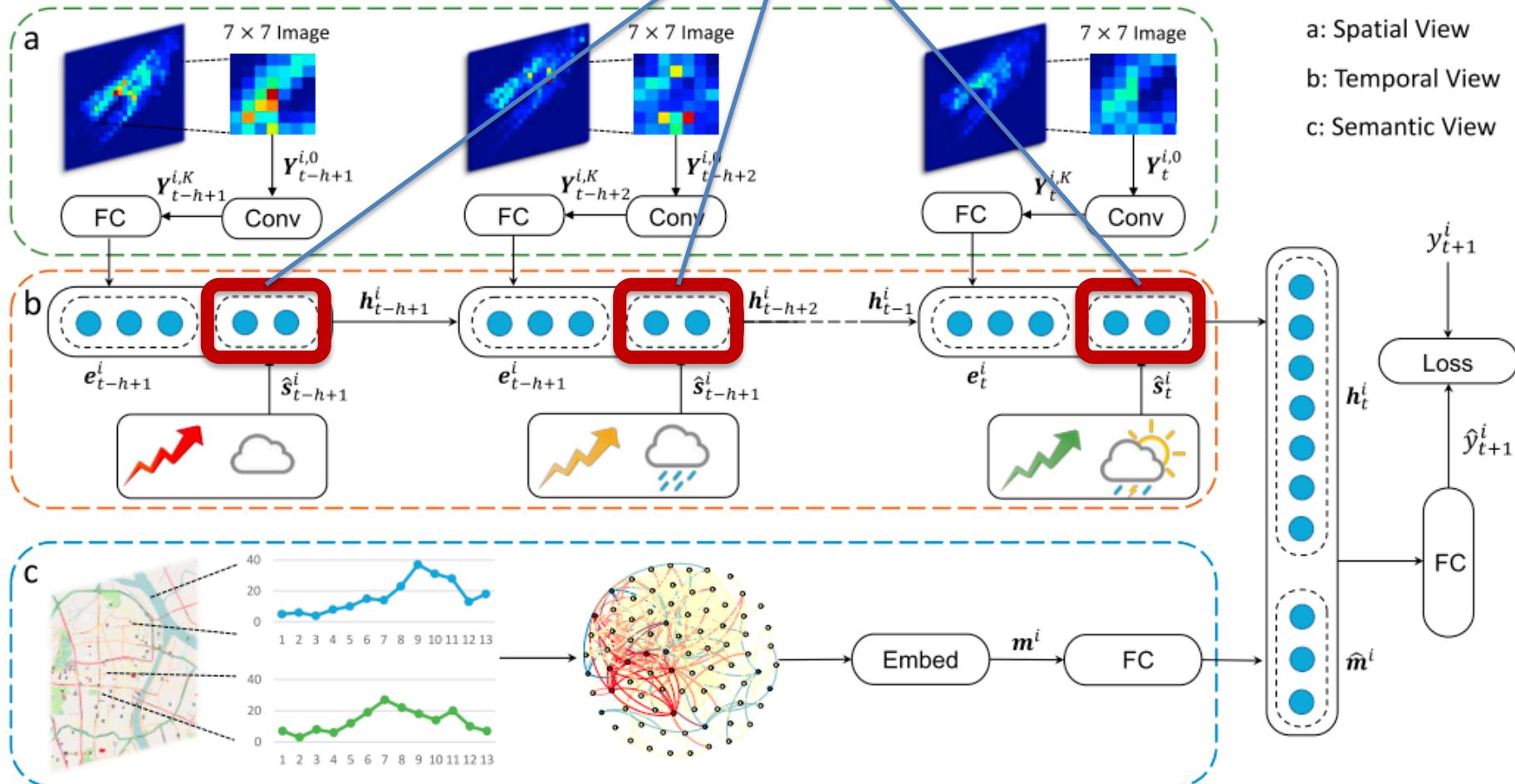
Pennsylvania State University
jessieli@ist.psu.edu

Taxi demand prediction



Taxi demand prediction

?



- a: Spatial View
- b: Temporal View
- c: Semantic View

Embeddings

- We often use categorical variables in our models
 - Weather:
 - Rain, snow, sunny, cloudy
 - Day of the week
 - Type of event
 - Concert, soccer game, basketball game
 - Service status
 - Normal, minor disruptions, major disruptions
- How to represent this type of data in our model?

Embeddings

- We often use categorical variables in our models
- Typically use **one-hot encoding**
 - Rain: [0, 1, 0, 0, 0, 0]
 - Snow [1, 0, 0, 0, 0, 0]
 - Sunny [0, 0, 1, 0, 0, 0]

Embeddings

- We often use categorical variables in our models
- Typically use **one-hot encoding**

- Rain:

[0, 1, 0, 0, 0, 0]

- Snow

[1, 0, 0, 0, 0, 0]

- Sunny

[0, 0, 1, 0, 0, 0]

Vocabulary size = 6

Embeddings

Issues with one-hot encoding

- For large vocabularies, it becomes very high dimensional
- Extremely sparse
 - DL methods don't work well with sparse data
- The mapping is completely uninformed: “similar” categories are not placed closer to each other in embedding space.

Embeddings

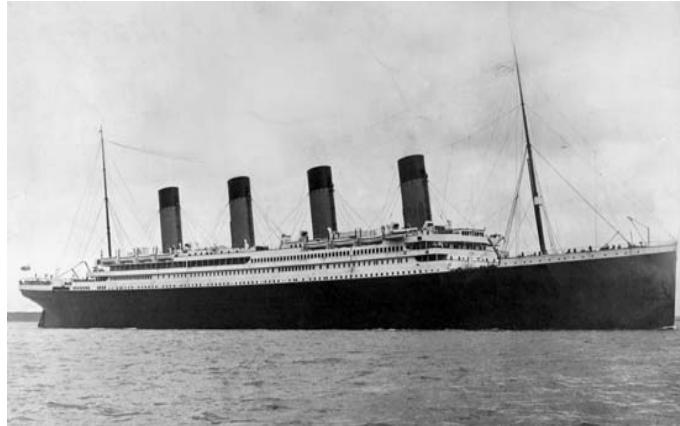
NLP tasks are faced with a similar problem

- How should we represent each word?
- Some words are more “similar”

Embeddings

NLP tasks are faced with a similar problem

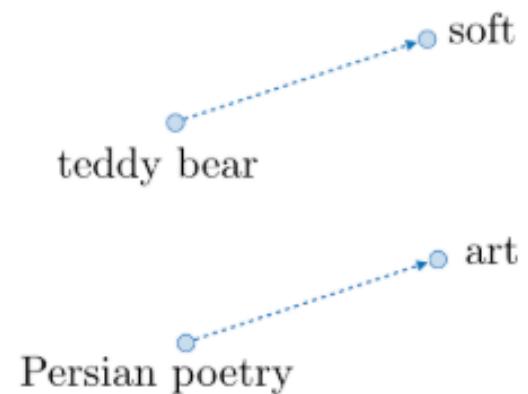
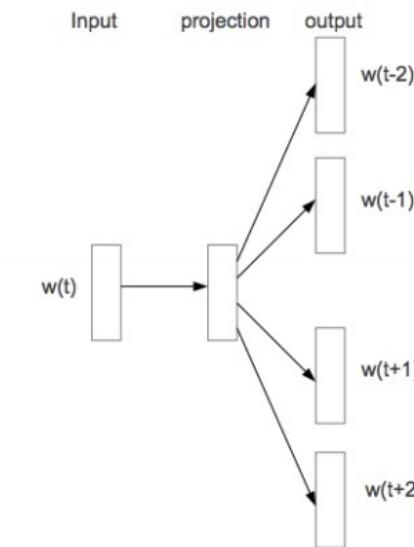
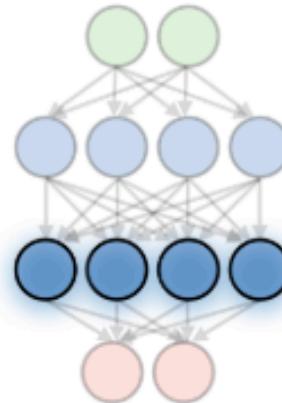
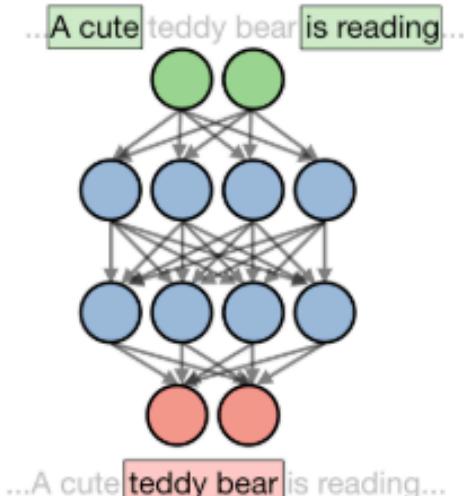
- How should we represent each word?



Embeddings

NLP tasks are faced with a similar problem

- Word2Vec



Train network on proxy task



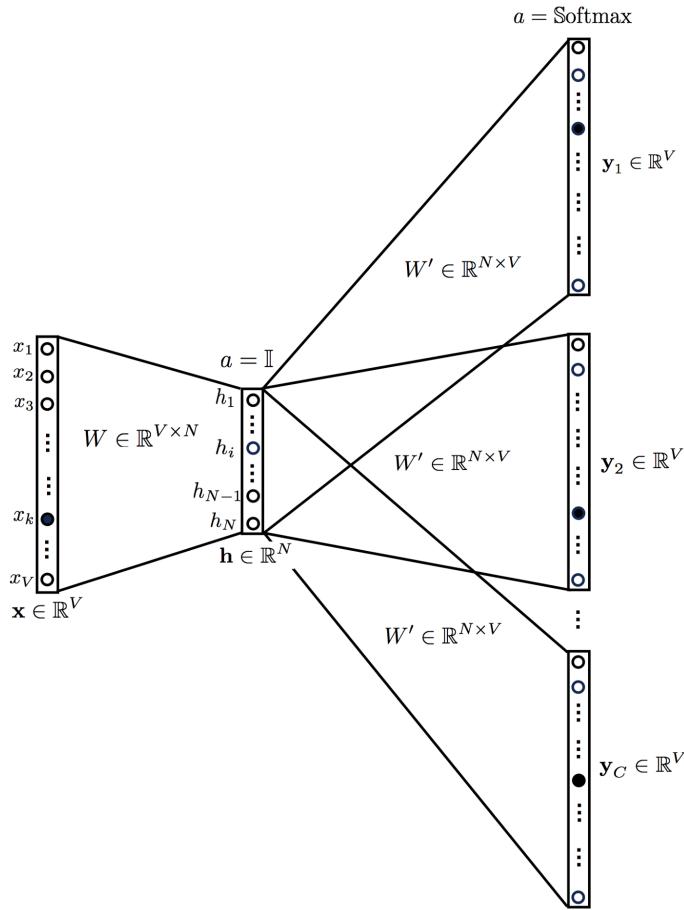
Extract high-level representation



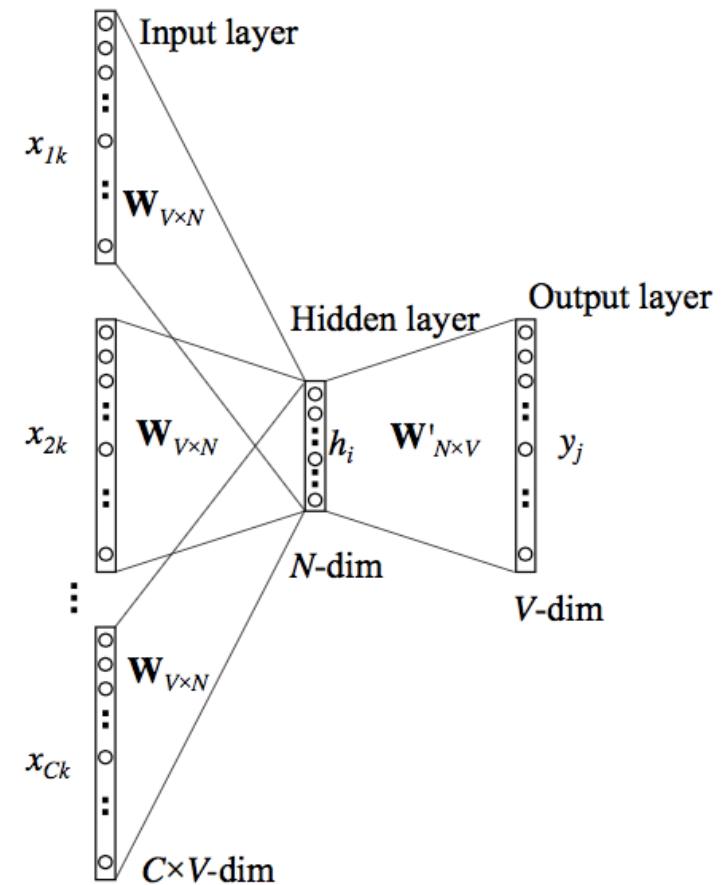
Compute word embeddings

<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks>

Skip Gram



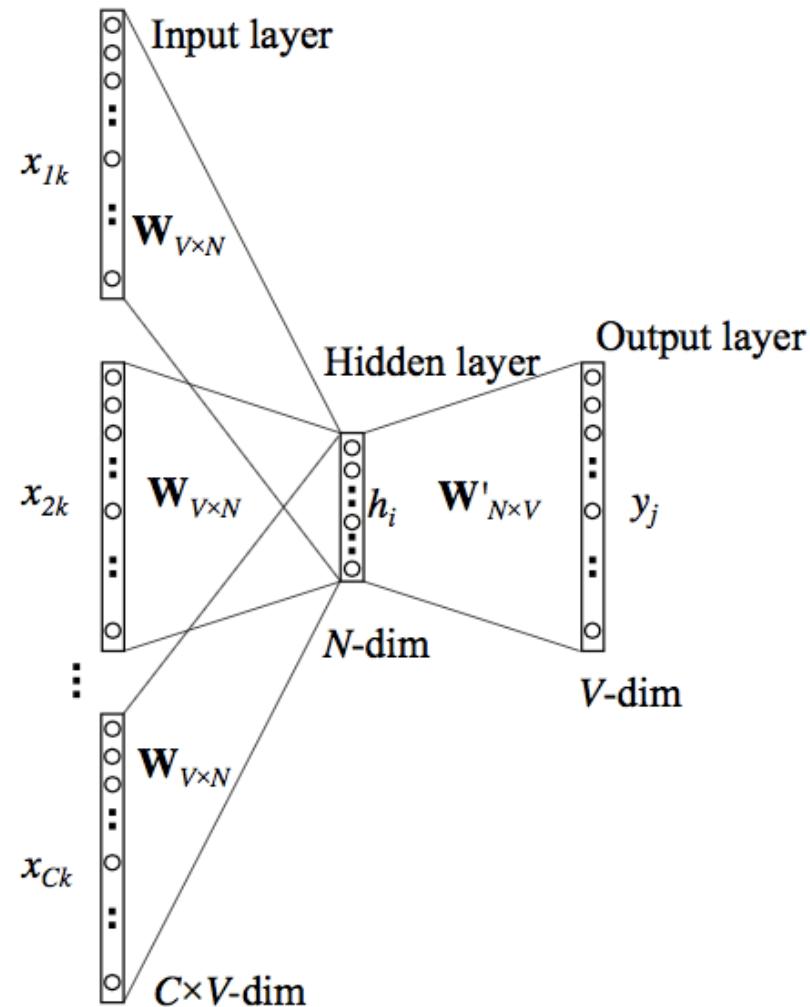
CBOW (Continuous Bag of Words)



Embeddings

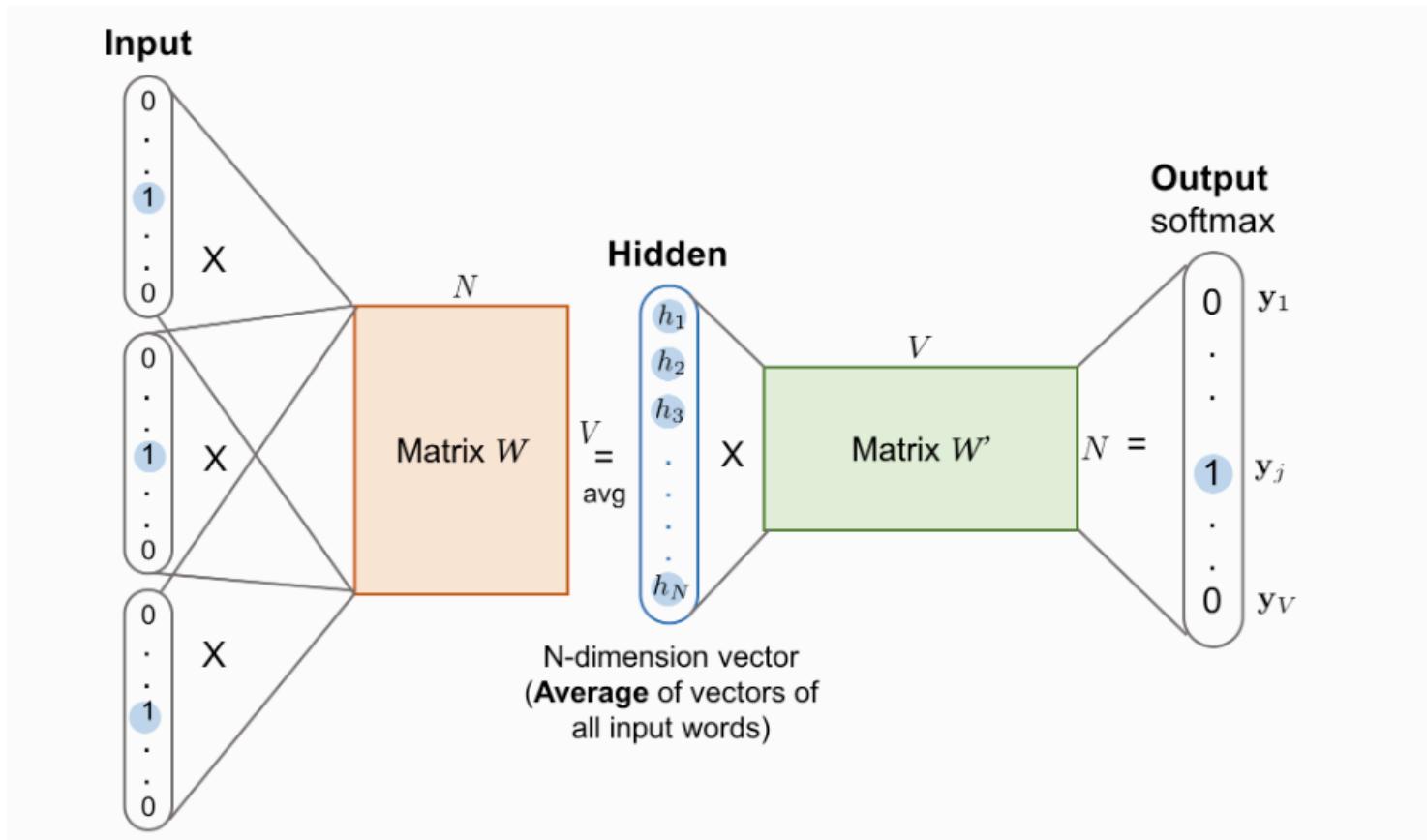
CBOW (Continuous Bag of Words)

- One-hot each word
- SoftMax layer at the end
- C: Context
- V: Vocab dimension
- N: Embedding size



Embeddings

$$\mathcal{L}_\theta = -\log \frac{\exp(v'_{w_O}^\top v_{w_I})}{\sum_{i=1}^V \exp(v'_{w_i}^\top v_{w_I})} = -v'_{w_O}^\top v_{w_I} + \log \sum_{i=1}^V \exp(v'_{w_i}^\top v_{w_I})$$



Embeddings

- Word2vec vs one-hot encoding

- Rain

- [0, 1, 0, 0, 0, 0]

- Rain

- [0.23, 0.051, 0.674, 0.599, 0.111, 0.672]

Embeddings

```
In [1]: import gensim
```

```
In [3]: model=gensim.models.Word2Vec.load("gensimModel")
```

```
In [4]: model.wv.most_similar(positive=['customer'])
```

```
Out[4]: [('passenger', 0.855195164680481),  
          ('lady', 0.6993716955184937),  
          ('person', 0.6990634202957153),  
          ('woman', 0.6973735690116882),  
          ('vagrant', 0.6603719592094421),  
          ('man', 0.6498541831970215),  
          ('child', 0.6353968977928162),  
          ('male', 0.6159298419952393),  
          ('passengers', 0.606838583946228),  
          ('female', 0.5832453966140747)]
```

```
In [14]: model.wv.most_similar(positive=['oxford'], negative=['incident'])
```

```
Out[14]: [('nhg', 0.484191358089447),  
           ('vic', 0.45740655064582825),  
           ('leicester', 0.4572640061378479),  
           ('n241', 0.4489314556121826),  
           ('pimlico', 0.4430849552154541),  
           ('s220', 0.4392928183078766),  
           ('stockwell', 0.43596023321151733),  
           ('s221', 0.4301430583000183),  
           ('s250', 0.42746976017951965),  
           ('n244', 0.4272845387458801)]
```

```
In [16]: model.wv.most_similar(positive=['incident'])
```

```
Out[16]: [('matter', 0.6234616041183472),  
           ('spad', 0.42893028259277344),  
           ('issue', 0.4024146795272827),  
           ('situation', 0.39529454708099365),  
           ('item', 0.3943479061126709),  
           ('conversation', 0.39270853996276855),  
           ('discussion', 0.39042913913726807),  
           ('sm', 0.38687431812286377),  
           ('soo', 0.3864406943321228),
```

Embeddings

- Word vectors demonstrated the usefulness of dense vector representation of the inputs
- In other domains, we simply initialize each categorical variable as a list of random numbers, then train it with other parameters in the model

Lab 1

Keras

TensorFlow / Theano / CNTK / ...

CUDA / cuDNN

BLAS, Eigen

GPU

CPU

Metro passenger flow prediction

- Predict passenger arrival at stations for the next 15 minutes



Contents lists available at [ScienceDirect](#)

Transportation Research Part C

journal homepage: www.elsevier.com/locate/trc



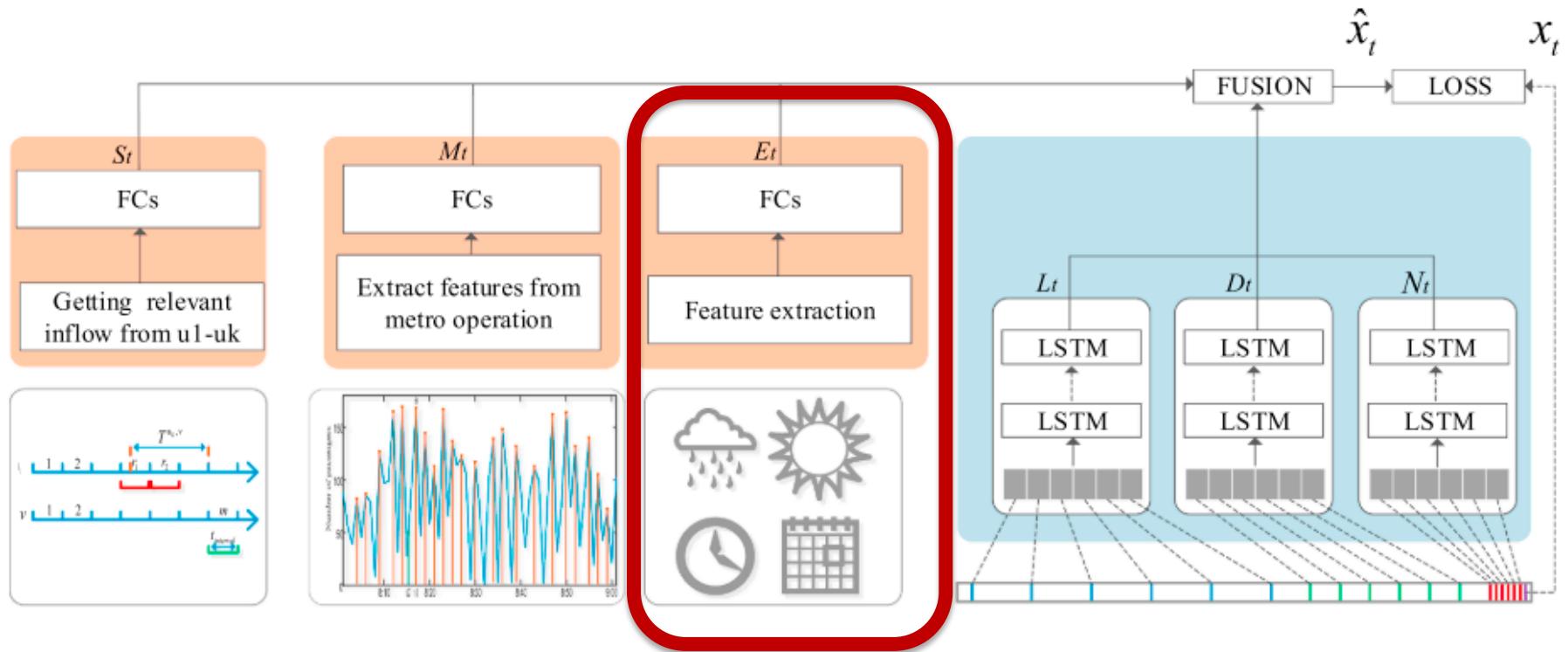
DeepPF: A deep learning based architecture for metro passenger flow prediction

Yang Liu, Zhiyuan Liu*, Ruo Jia



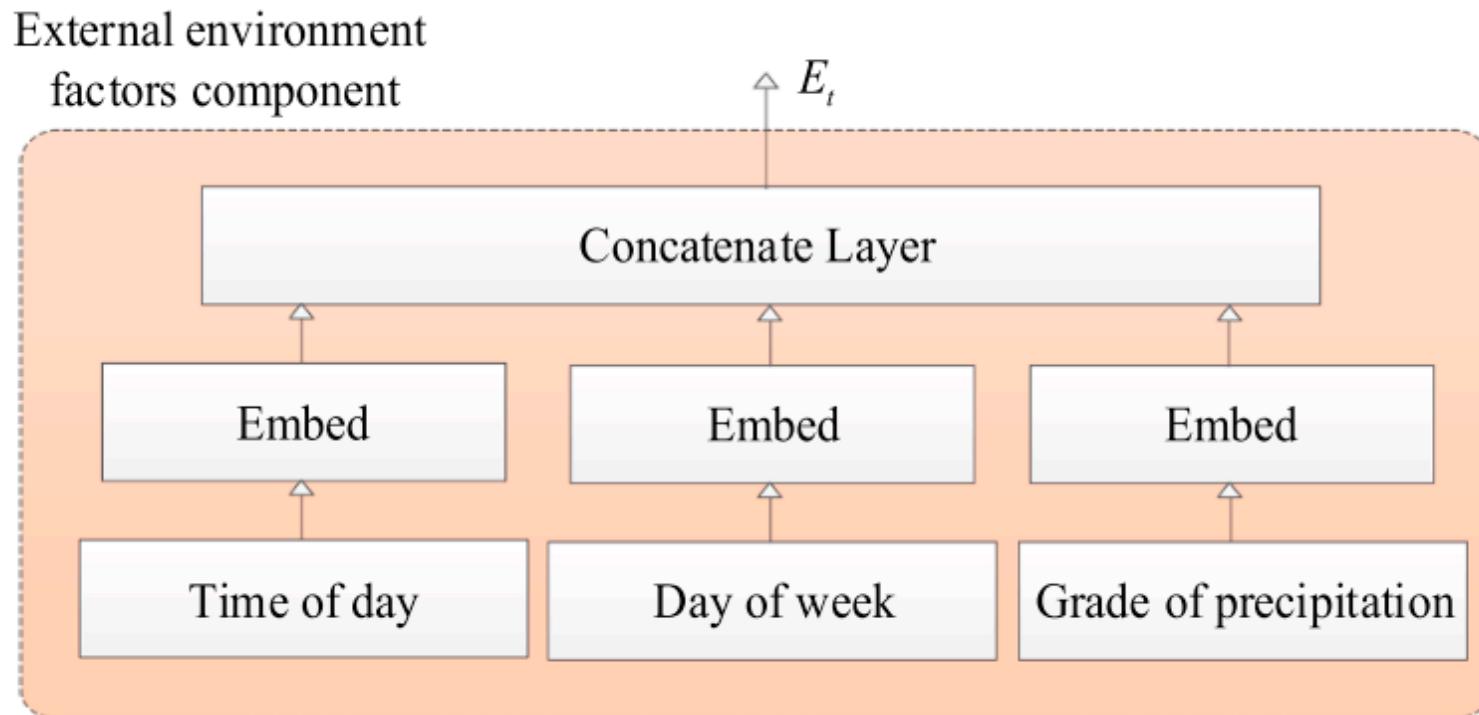
Metro passenger flow prediction

- Predict passenger arrival at stations for the next 15 minutes



Metro passenger flow prediction

- Predict passenger arrival at stations for the next 15 minutes



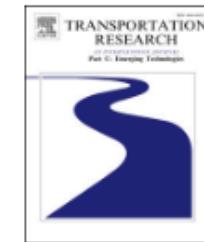
Demand prediction for bike sharing services



Contents lists available at [ScienceDirect](#)

Transportation Research Part C

journal homepage: www.elsevier.com/locate/trc



The station-free sharing bike demand forecasting with a deep learning approach and large-scale datasets



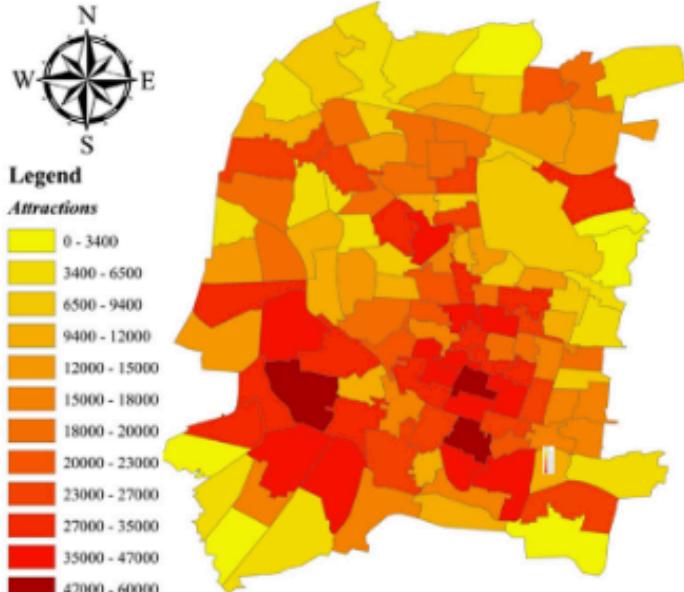
Chengcheng Xu*, Junyi Ji, Pan Liu



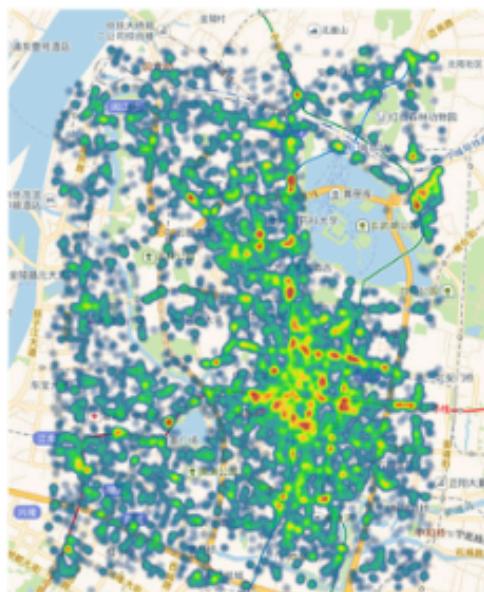
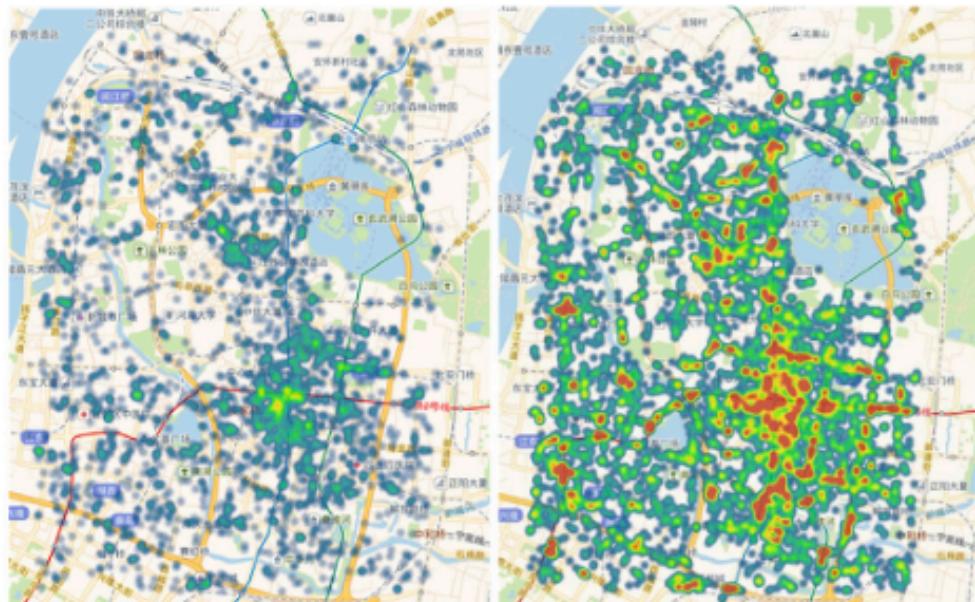
(a) Station-free bike sharing



(b) Bike sharing system with docking station



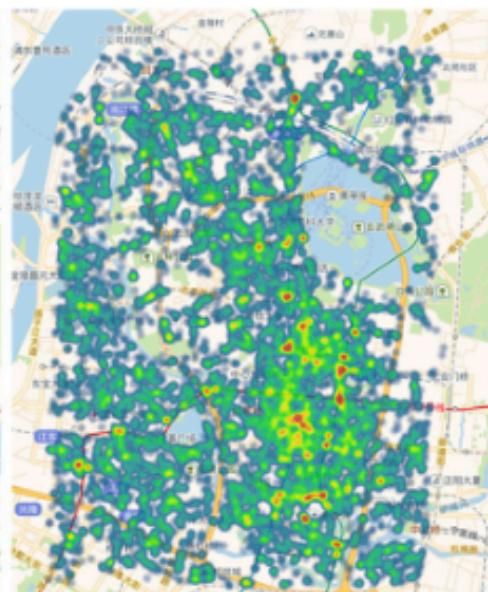
(a) Spatial distribution of Attraction



(d) Attraction 10:00-15:00



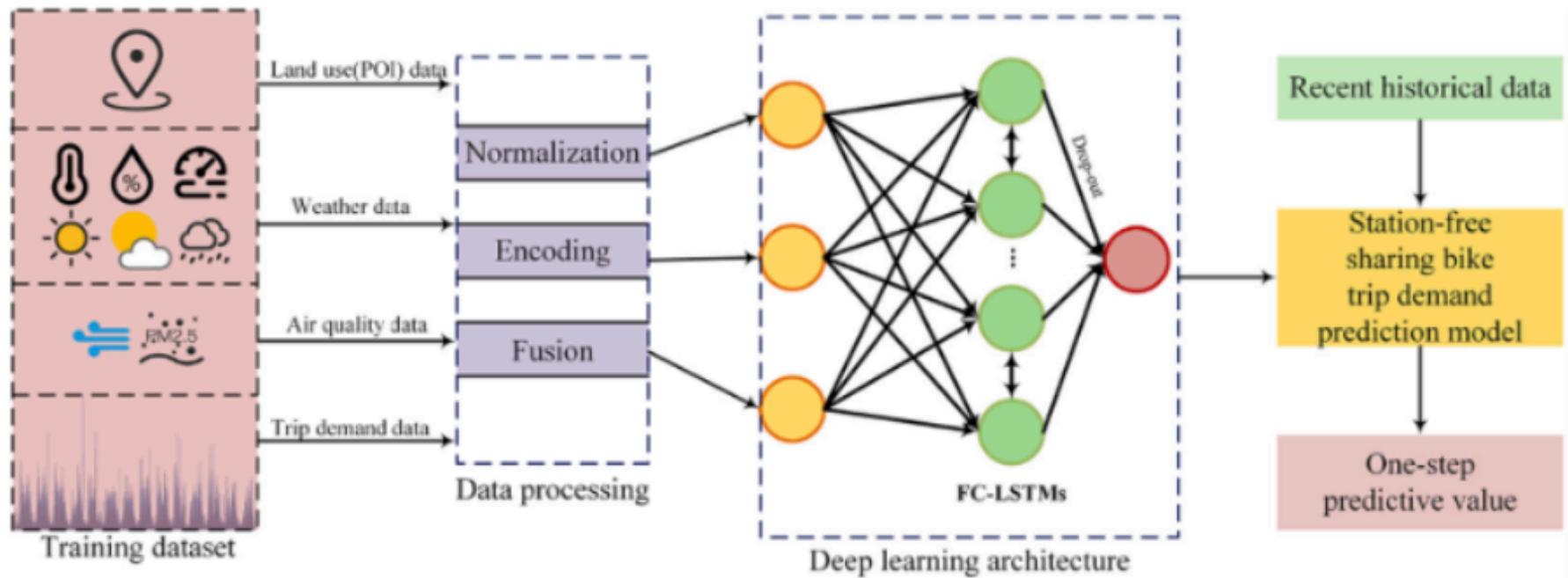
(e) Attraction 15:00-20:00

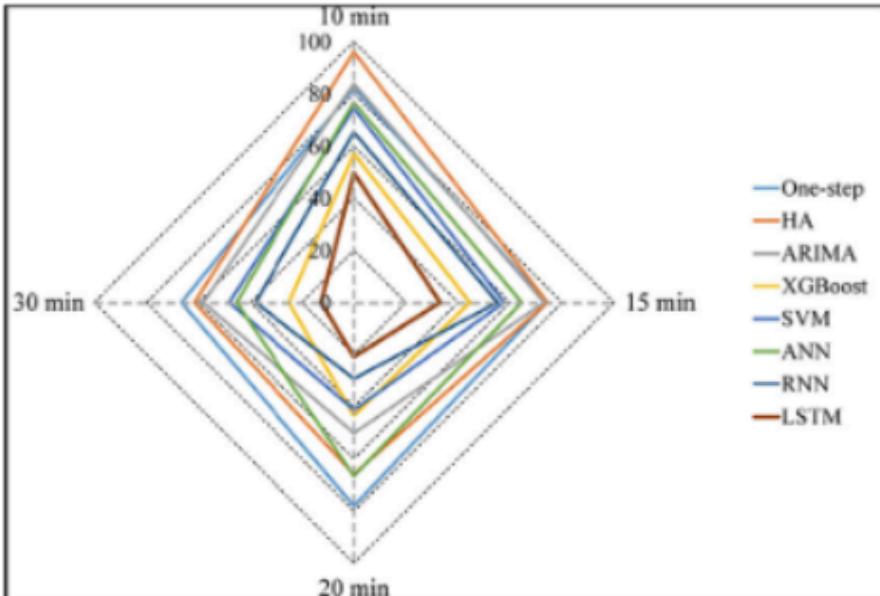


(f) Attraction 20:00-24:00

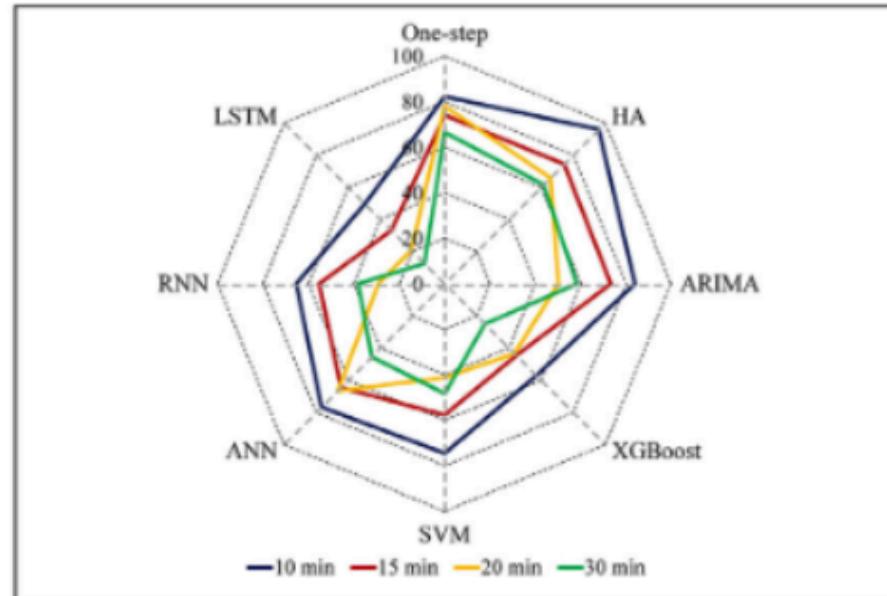
Demand prediction for bike sharing services

- Embedding + LSTM

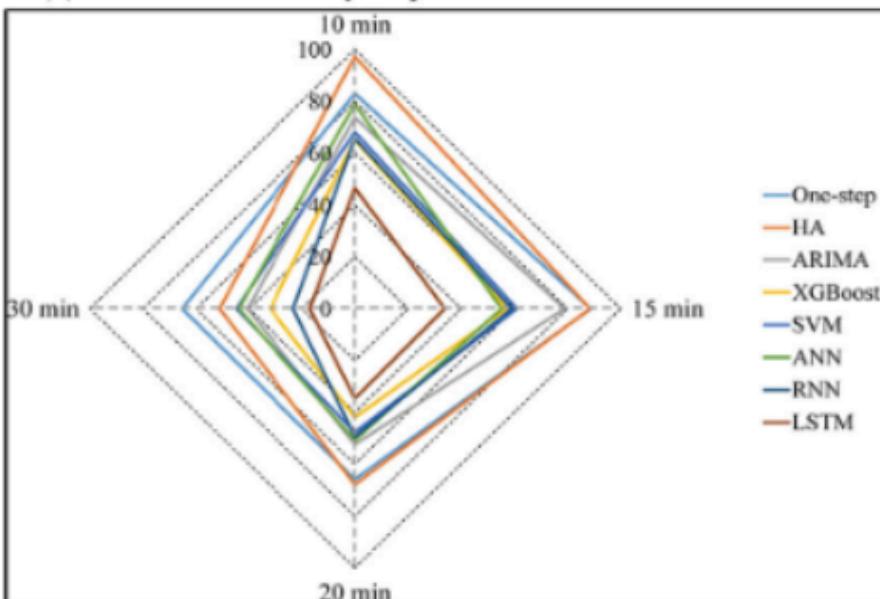




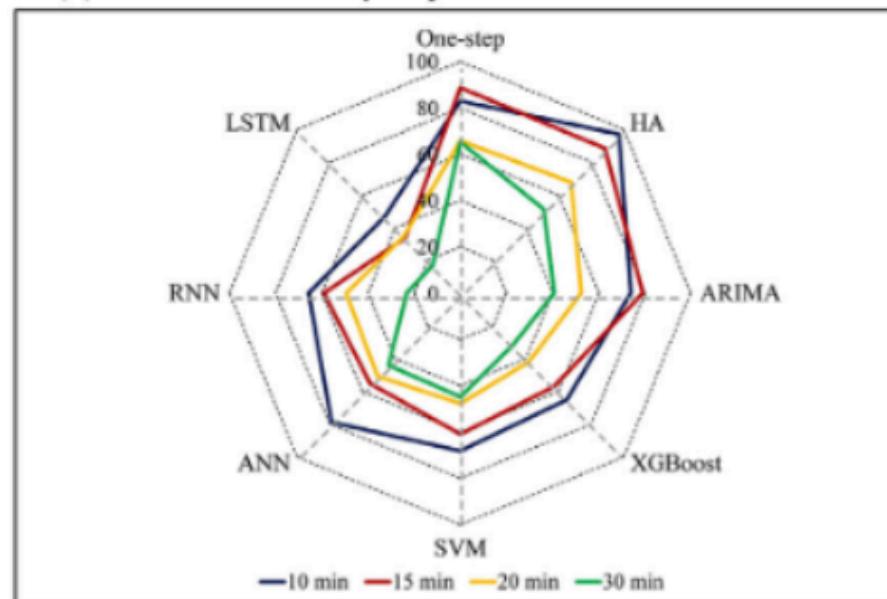
(a) Prediction accuracy of production at different intervals



(b) Prediction accuracy of production of different models



(c) Prediction accuracy of attraction at different intervals



(d) Prediction accuracy of attraction of different models

Mobility Embeddings

Unsupervised Learning of Parsimonious General-Purpose Embeddings for User and Location Modelling

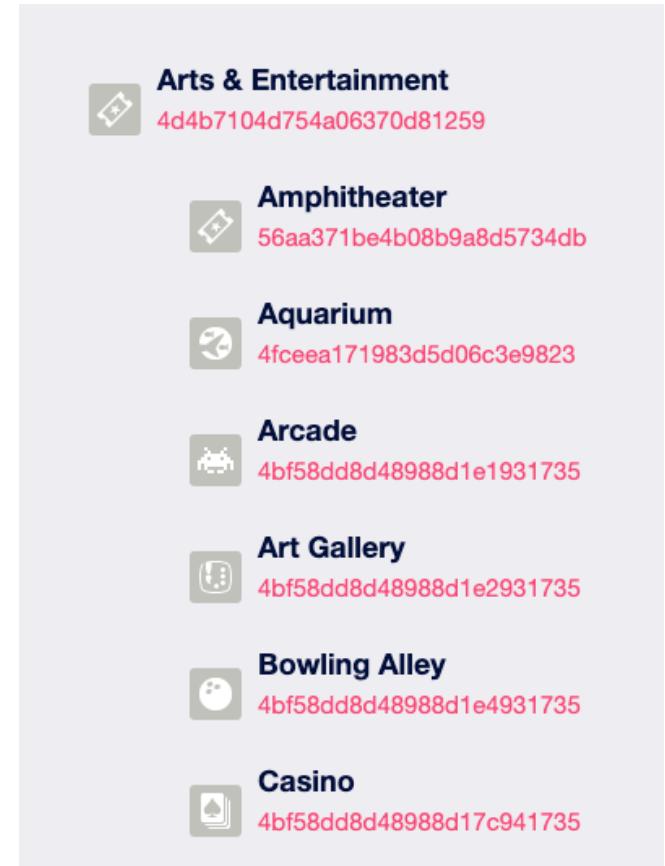
JING YANG, ETH Zurich, Switzerland

CARSTEN EICKHOFF, ETH Zurich, Switzerland

- Check-in data from Foursquare
- Learn
 - Geographic, temporal, and functional embeddings
- Application
 - Location recommendation
 - Urban functional zone study
 - Crime prediction

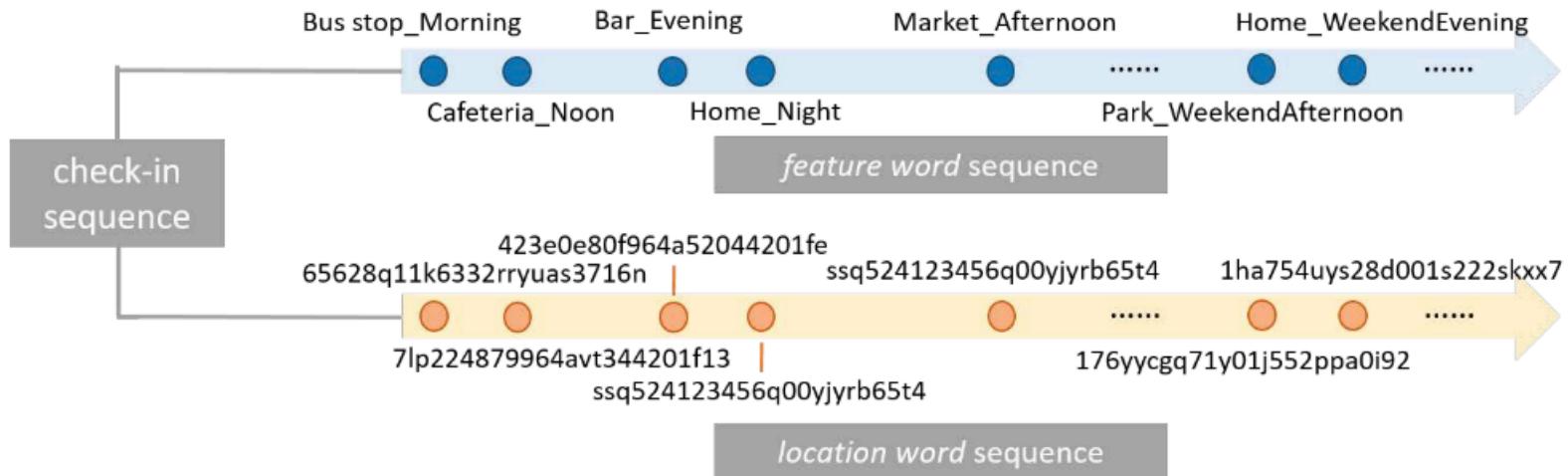
Mobility Embeddings

- Check-in is defines as a tuple $c = \langle u, f, t, l \rangle$
 - User u
 - Functional role f
 - Restaurant, bar, café, etc.
 - Time
 - Location (ID)



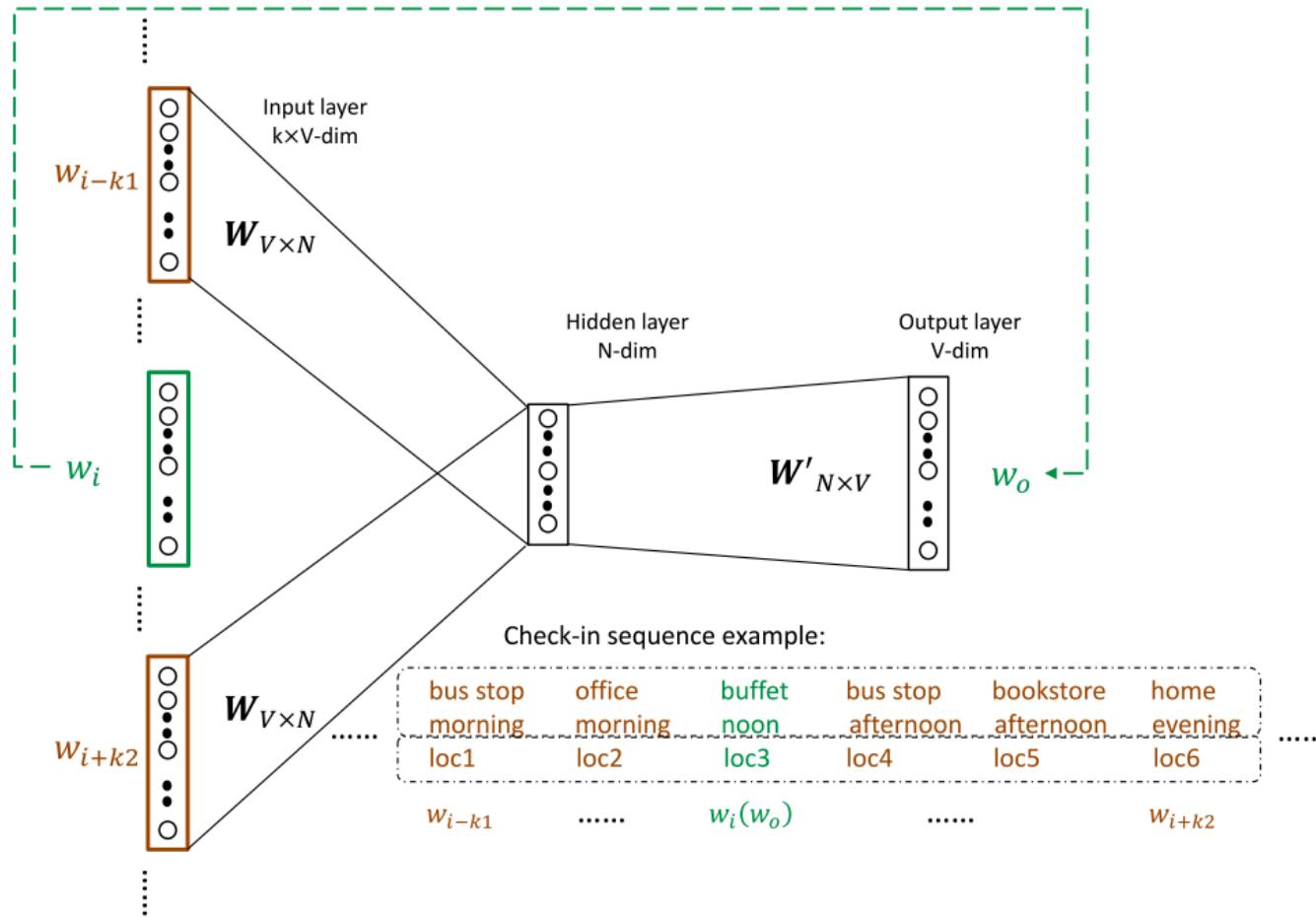
Mobility Embeddings

- Feature-word and Location-word sequences

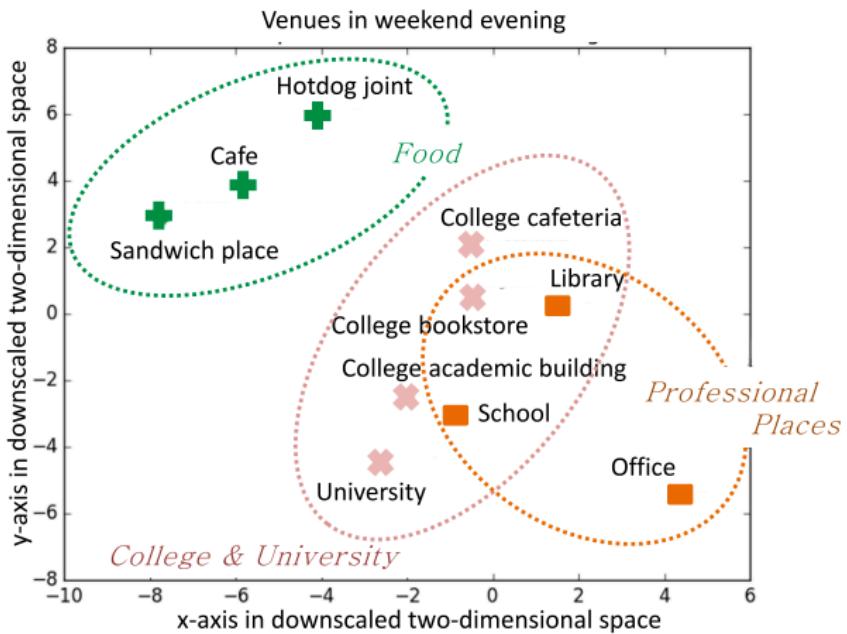


Mobility Embeddings

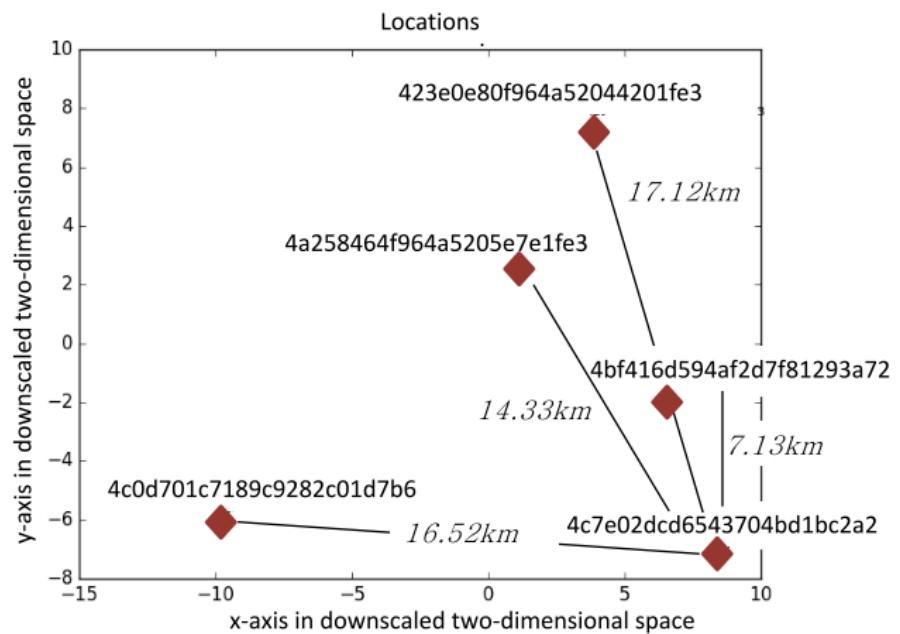
- Feature-word and Location-word sequences



Mobility Embeddings



(a) *feature word* embeddings



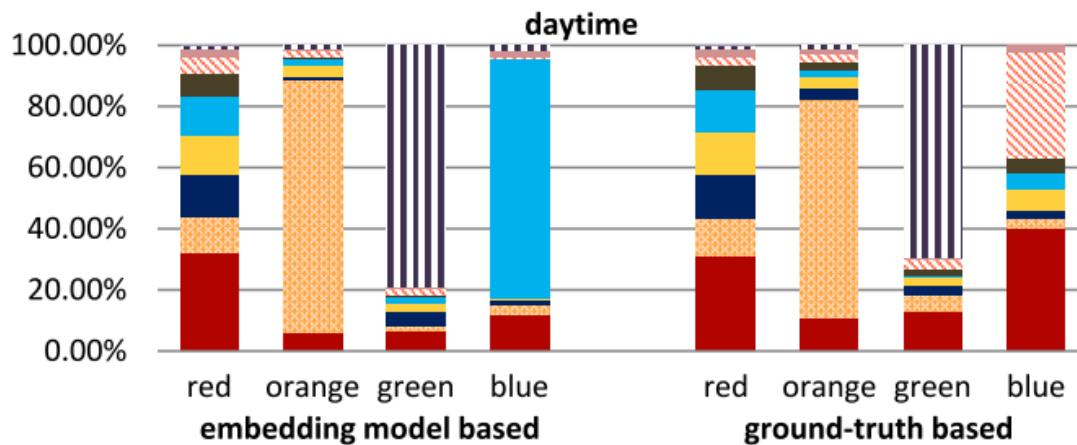
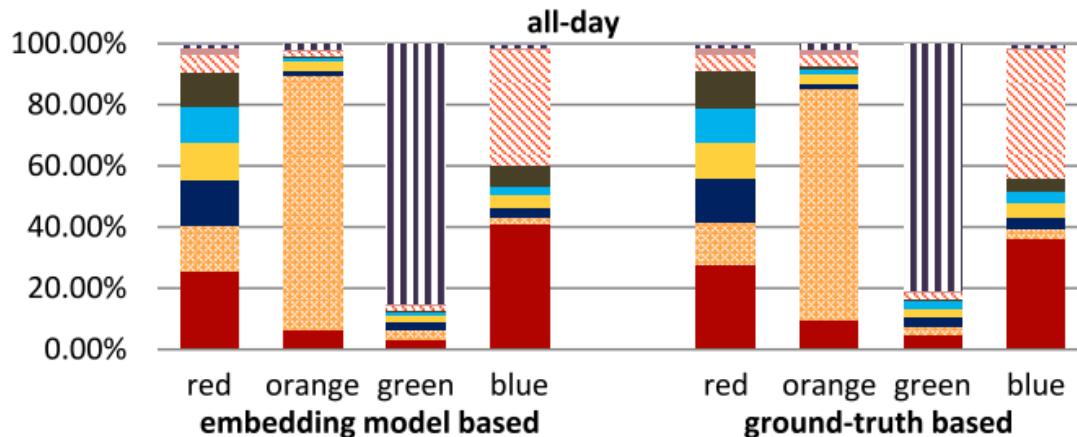
(b) *location word* embeddings

Mobility Embeddings

- Urban functional zone study
- Partition zones bases on the inhabitants' interactions with urban spaces
- Use feature-word embeddings with second-level venue categories
- Represent each neighborhood using the mean of all contained check-in vectors
- Feed into KNN



Mobility Embeddings



- Food
- Travel&Transport
- Outdoors&Recreation
- Shop&Service
- Professional Activity
- Arts&Entertainment
- Nightlife
- College&University
- Residence
- Event

Mobility Embeddings

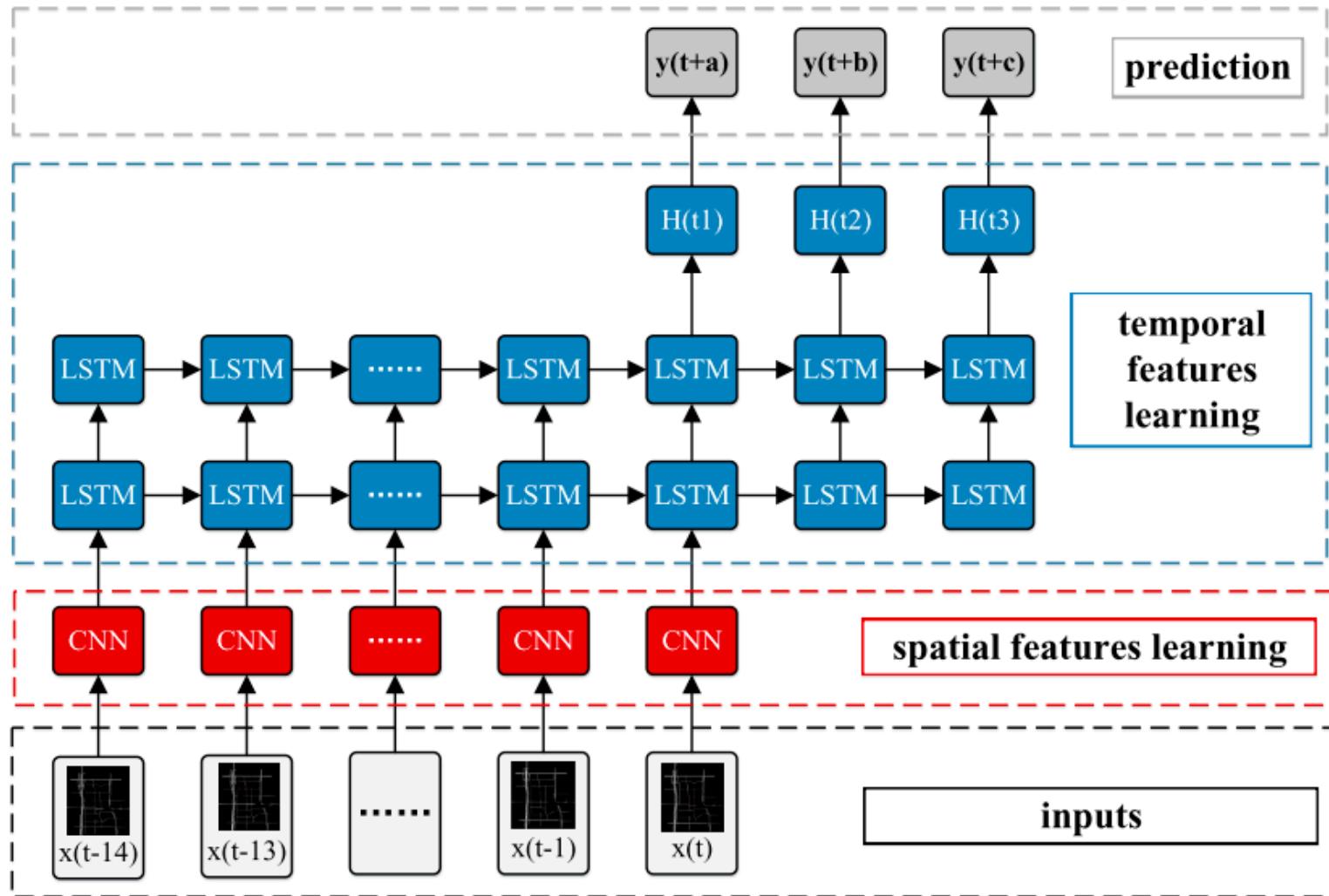


(a) neighborhood in the day

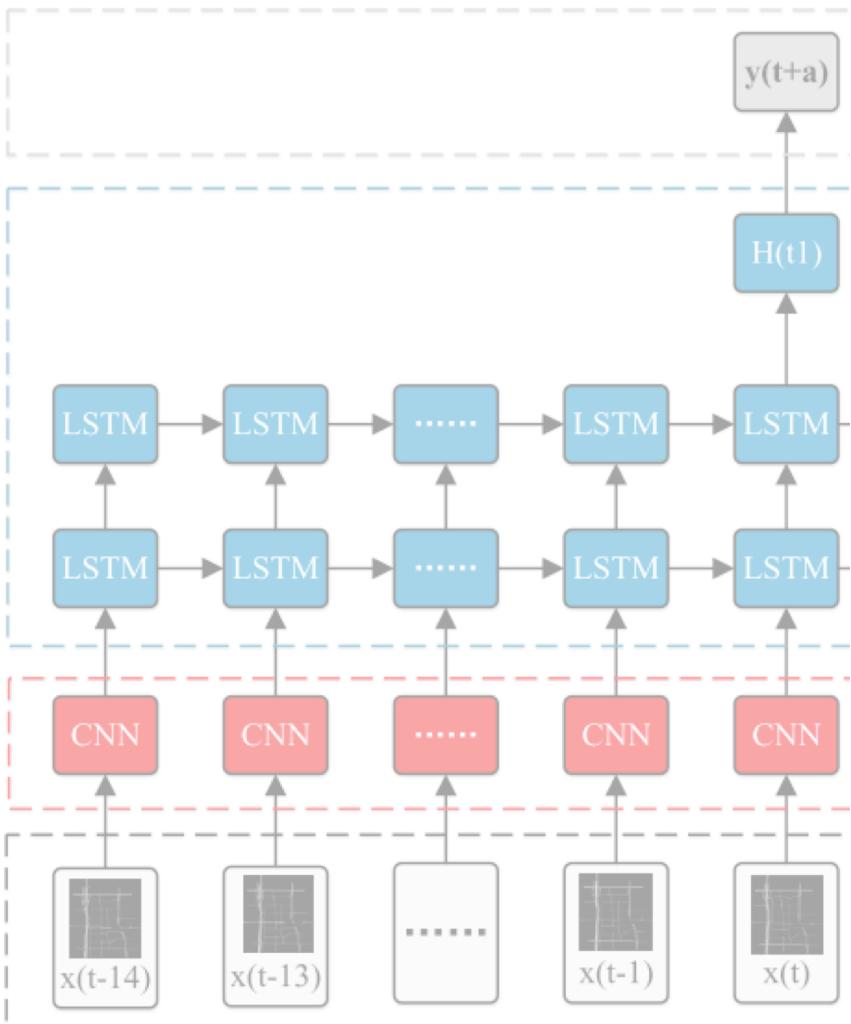


(b) neighborhood in the night

Two issues



Two issues

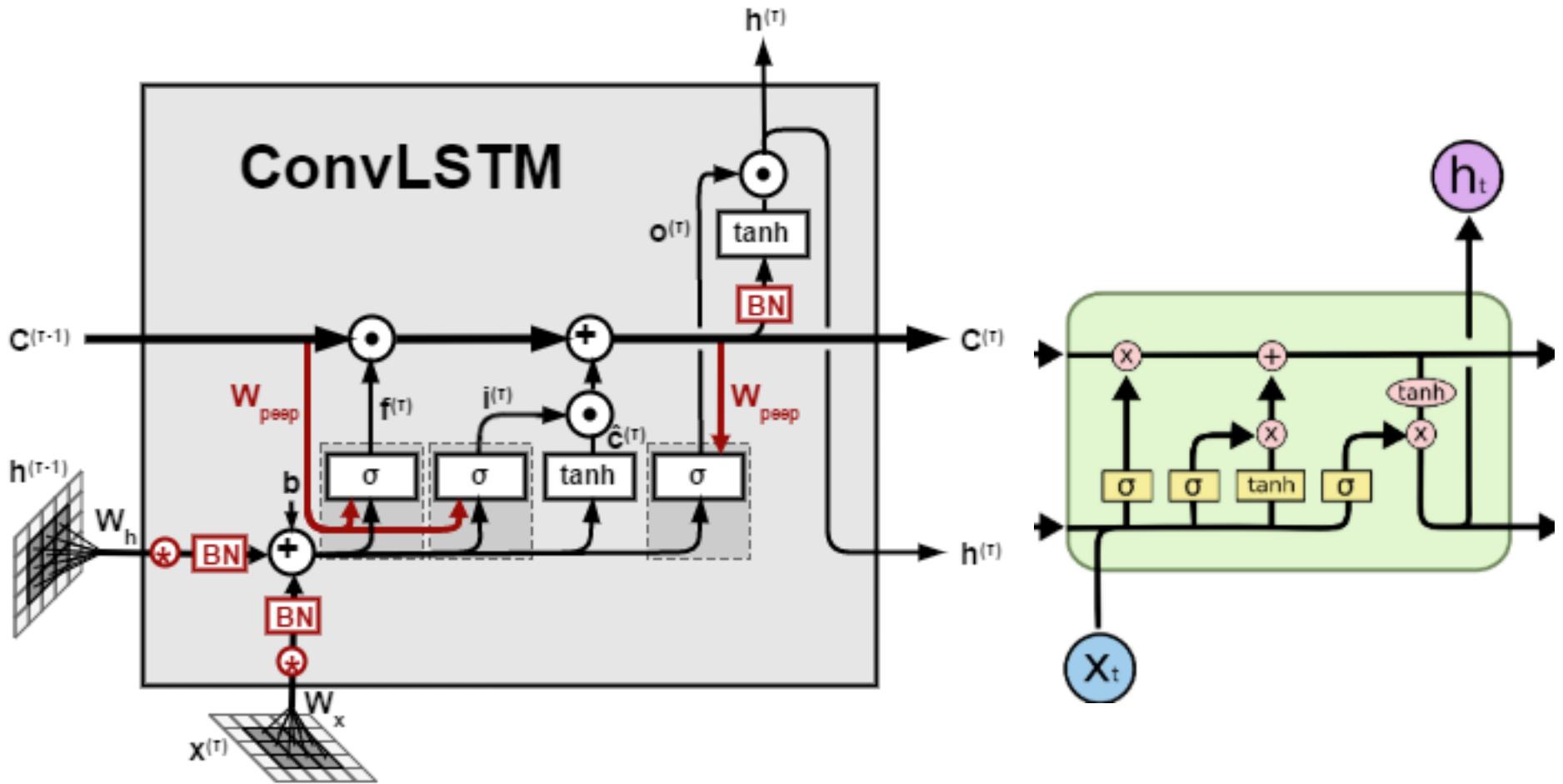


	0	Inputs	1	(163,148)
	Convolution	16		(3,3)
1	Max-pooling	16		(2,2)
	Activation (relu)	--		--
	Batch-normalization	--		--
2	Convolution	32		(3,3)
	Max-pooling	32		(2,2)
	Activation (relu)	--		--
	Batch-normalization	--		--
3	Convolution	64		(3,3)
	Activation (relu)	--		--
	Batch-normalization	--		--
4	Convolution	64		(3,3)
	Activation (relu)	--		--
	Batch-normalization	--		--
5	Convolution	128		(3,3)
	Max-pooling	128		(2,2)
	Activation (relu)	--		--
	Batch-normalization	--		--
6	Flatten			--
7	Fully connected			278
8	Lstm1			800
	Activation (tanh)			--
9	Lstm2			800
	Activation (tanh)			--
10	Dropout (0.2)			--
11	Fully connected			278

ConvLSTM

- The downside to using LSTMs for capturing the temporal dependencies is their inability in using the spatial information encoded in the input
- ConvLSTMs address this problem by having convolutional structures in both the input-to-state and state-to-state transitions.

ConvLSTM



ConvLSTM

- The downside to using LSTMs for capturing the temporal dependencies is their inability in using the spatial information encoded in the input
- ConvLSTMs address this problem by having convolutional structures in both the input-to-state and state-to-state transitions.

$$i_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \odot C_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \odot C_{t-1} + b_f)$$

$$C_t = f_{t-1} + i_t \odot \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + b_c)$$

$$o_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \odot C_t + b_o)$$

$$H_t = o_t \odot \tanh(C_t)$$

ConvLSTM

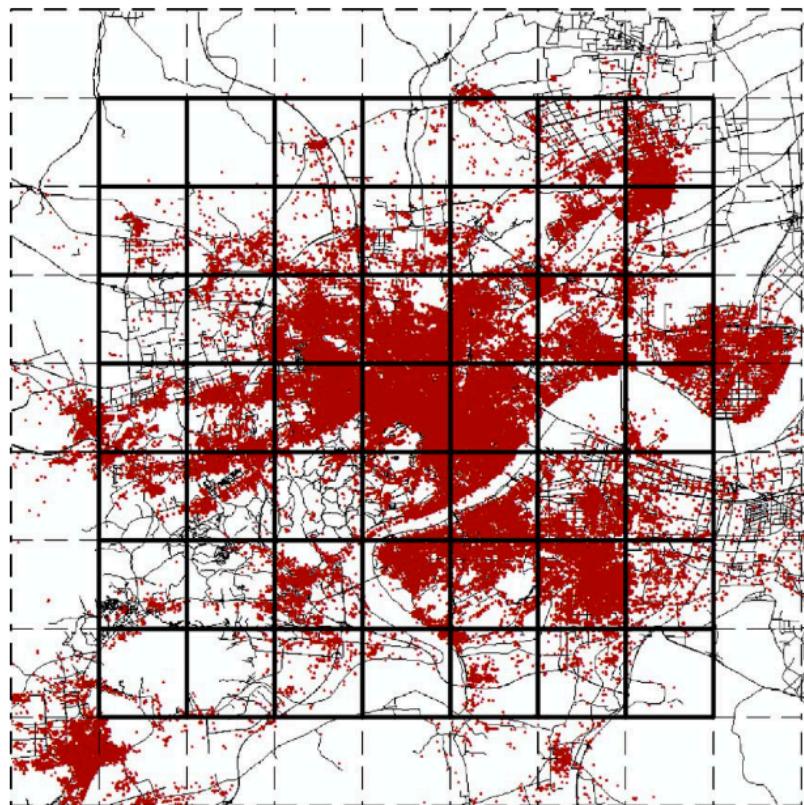
$$\begin{aligned} i_t &= \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \odot C_{t-1} + b_i) \\ f_t &= \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \odot C_{t-1} + b_f) \\ C_t &= f_t c_{t-1} + i_t \odot \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + b_c) \\ o_t &= \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \odot C_t + b_o) \\ H_t &= o_t \odot \tanh(C_t) \end{aligned}$$

LSTM

$$\begin{aligned} i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\ c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\ h_t &= o_t \tanh(c_t) \end{aligned}$$

ConvLSTM

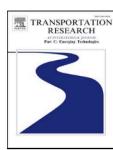
- Demand prediction for ride hailing services
- 1 million requests btw 2015-2016
- 1-hour time intervals
- $4.77 * 4.88$ km zones
- Temperature, humidity, weather state, wind speed, and visibility



Contents lists available at [ScienceDirect](#)

Transportation Research Part C

journal homepage: www.elsevier.com/locate/trc



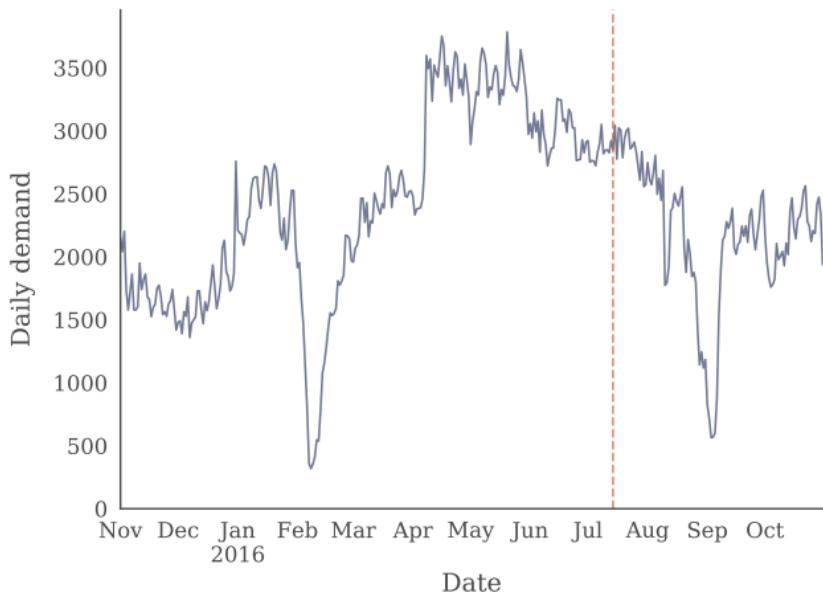
Short-term forecasting of passenger demand under on-demand ride services: A spatio-temporal deep learning approach



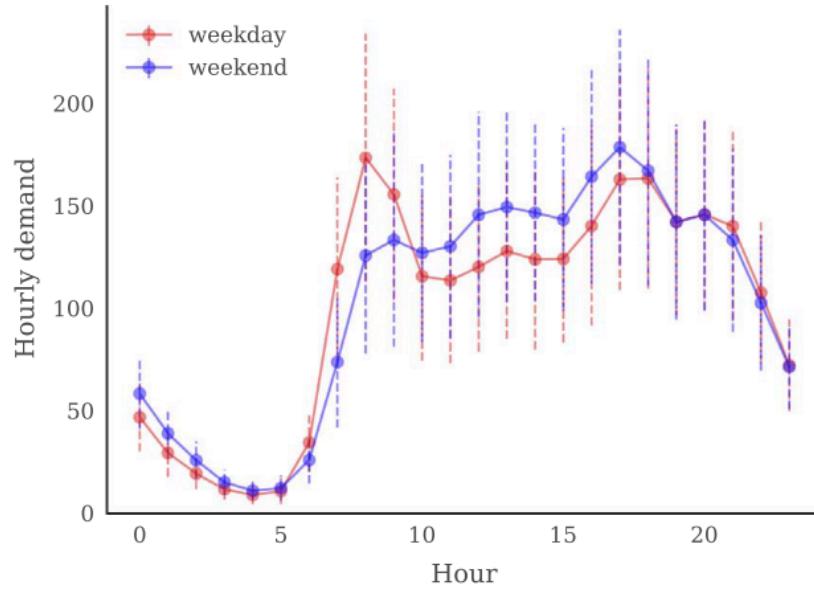
Jintao Ke^a, Hongyu Zheng^b, Hai Yang^a, Xiqun (Michael) Chen^{b,*}

ConvLSTM

- Demand prediction for ride hailing services

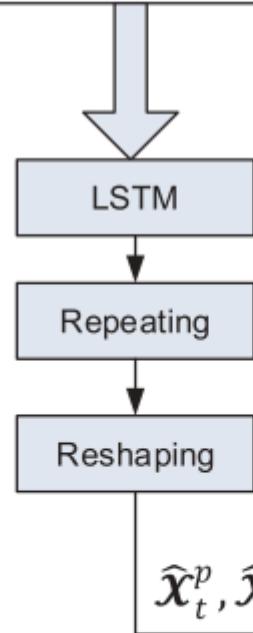
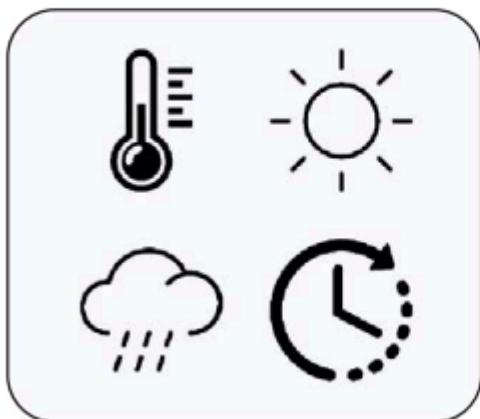


(a) Demand in different days

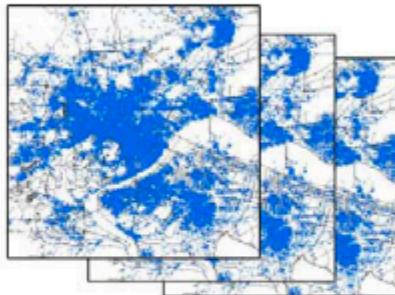


(b) Mean and standard deviation of hourly demand

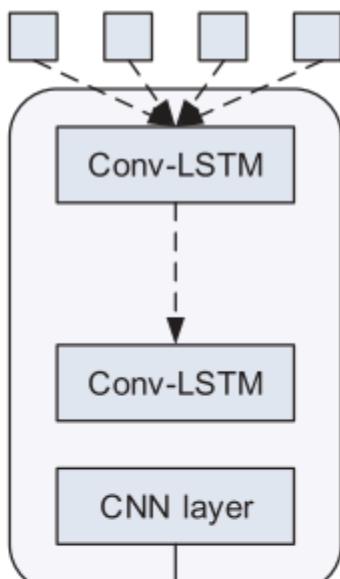
Sequences of time and weather data



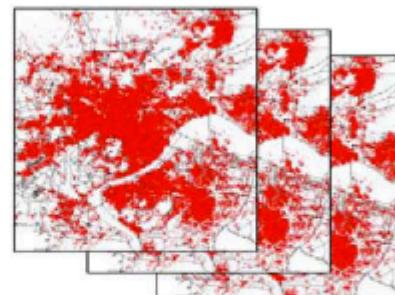
Map sequences of travel time rate



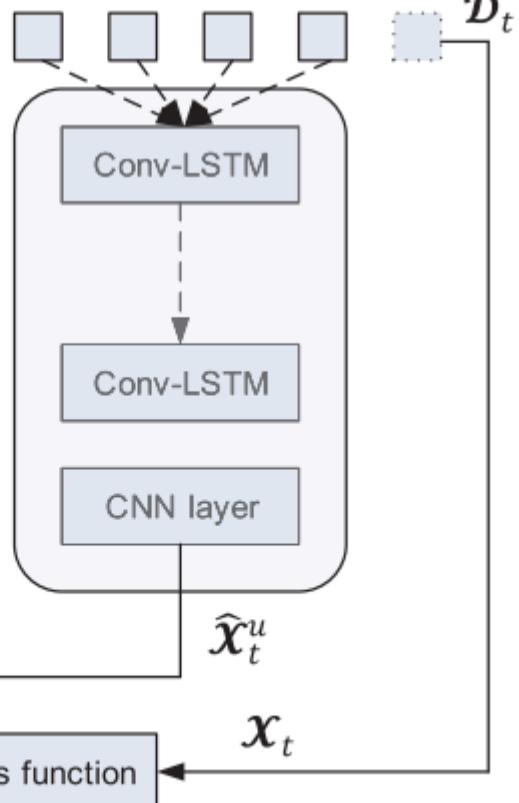
$$\mathcal{T}_{t-K_t}, \mathcal{T}_{t-K_t+1}, \dots, \mathcal{T}_{t-1}$$



Map sequences of demand intensity



$$\mathcal{D}_{t-K_d}, \mathcal{D}_{t-K_d+1}, \dots, \mathcal{D}_{t-1}$$



Fusion



\hat{x}_t

Loss function

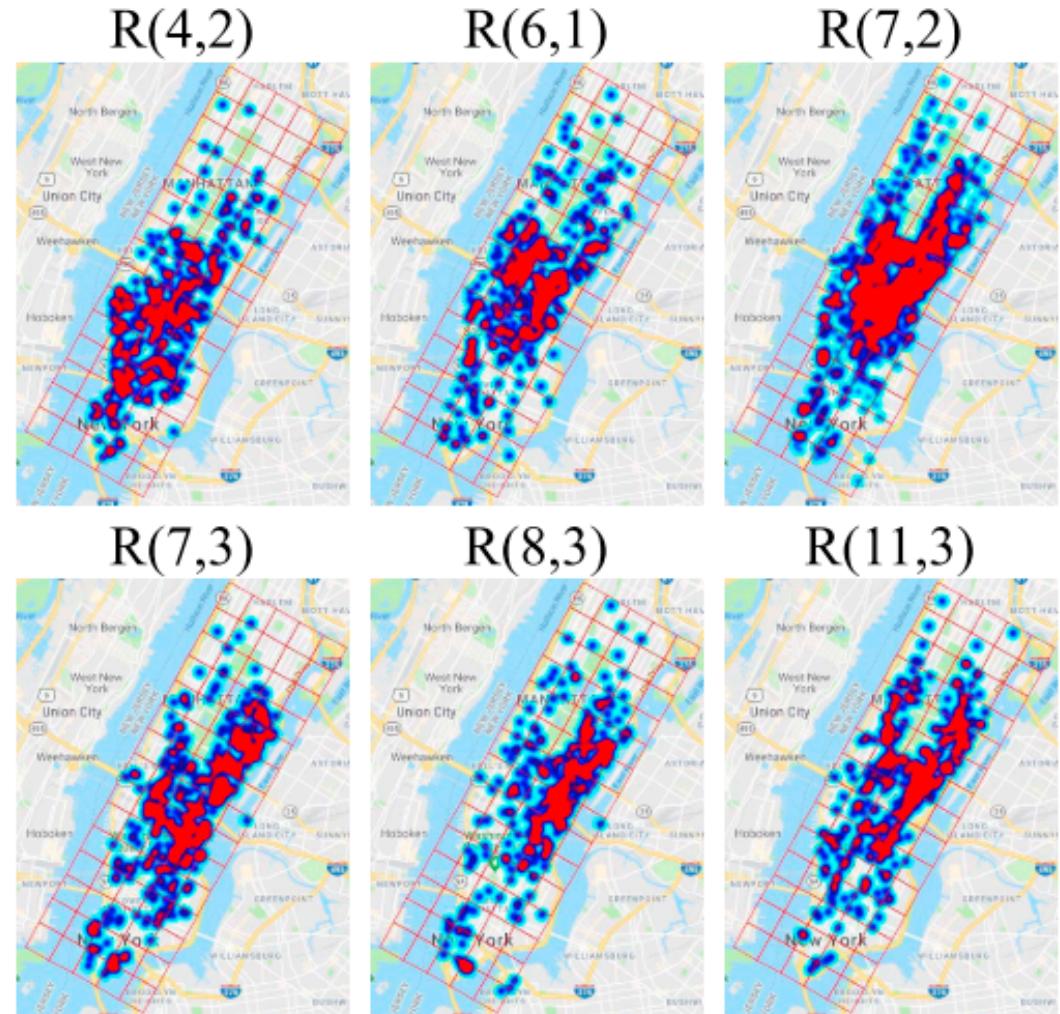
x_t

Contextualized Spatial–Temporal Network for Taxi Origin-Destination Demand Prediction

Lingbo Liu^{ID}, Zhilin Qiu, Guanbin Li^{ID}, *Member, IEEE*, Qing Wang, Wanli Ouyang^{ID}, *Senior Member, IEEE*, and Liang Lin^{ID}, *Senior Member, IEEE*

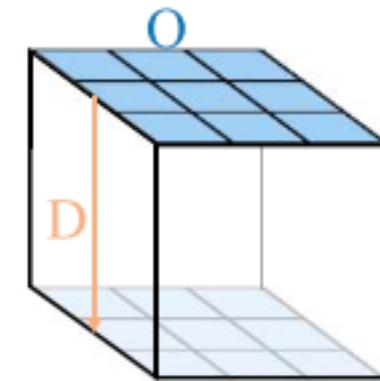
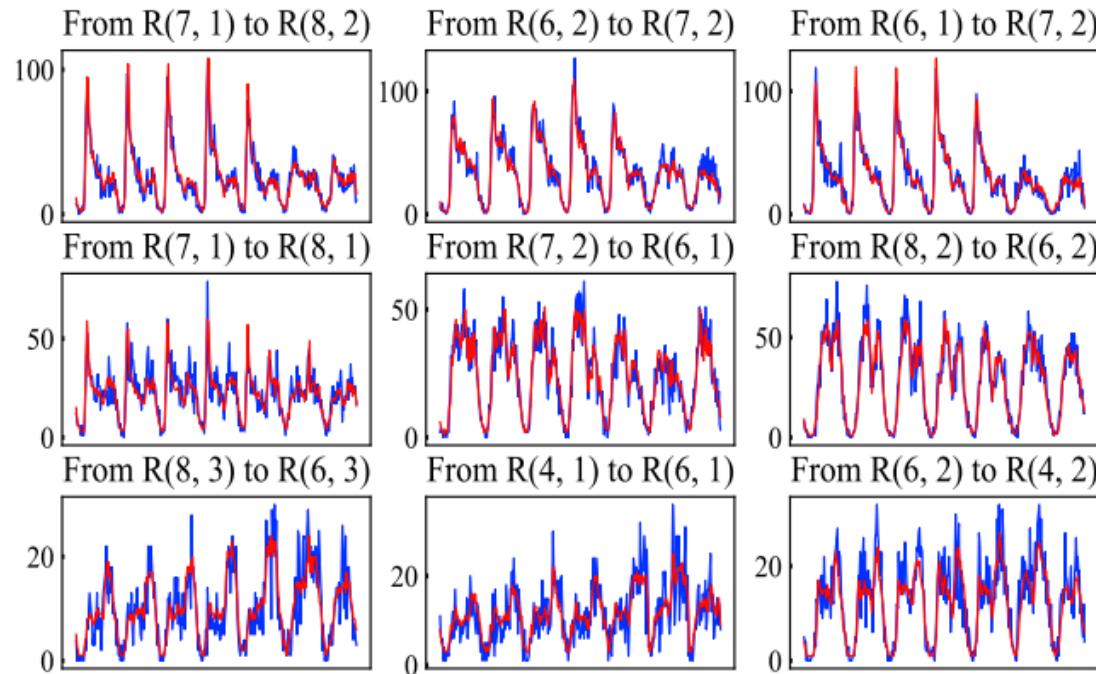
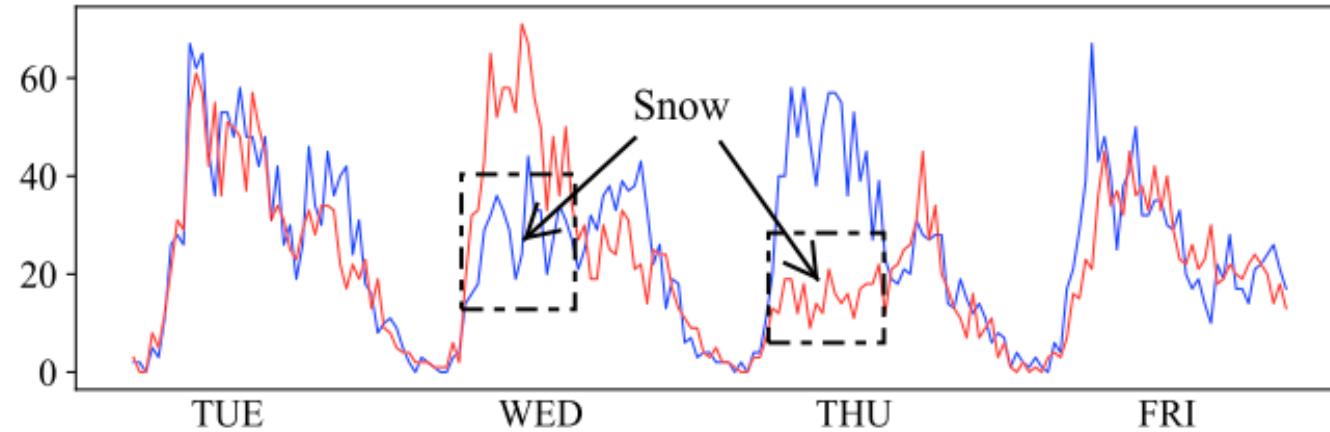


(a)

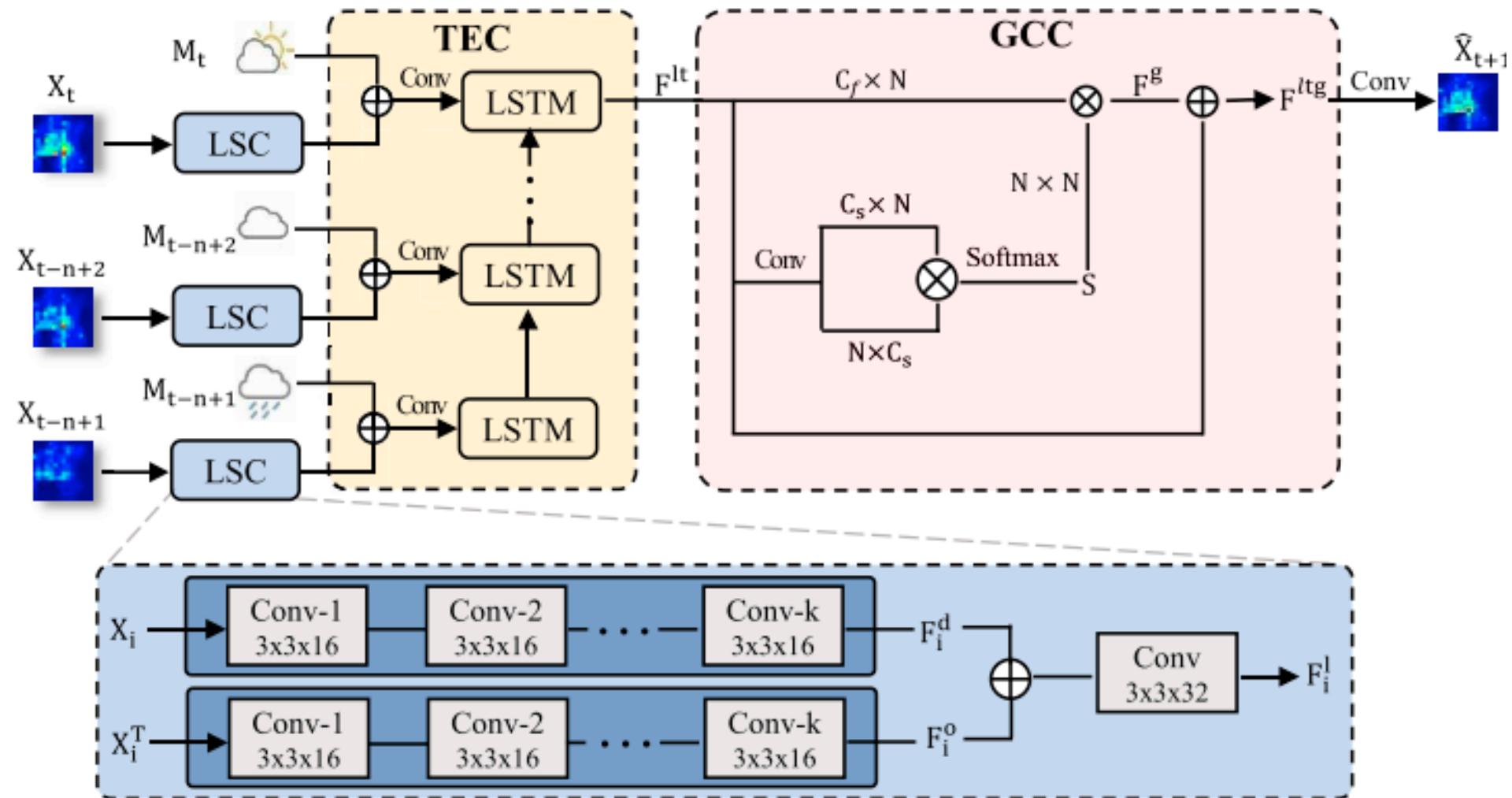


(b)

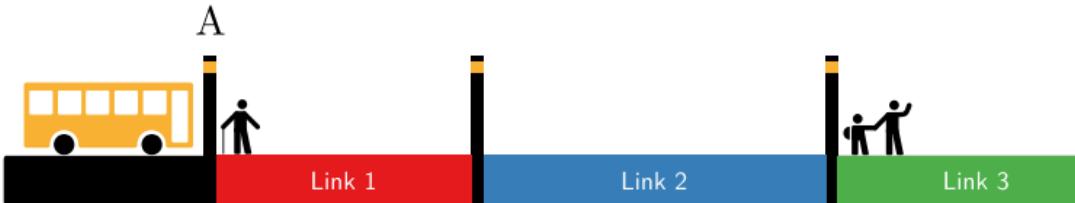
Demand prediction



Demand prediction



Bus travel time prediction



- AVL data from buses
- 1,2M travel time observations for the “4A” bus line in the period May to October 2017.
- aggregated into 15-min intervals



Multi-output bus travel time prediction with convolutional LSTM neural network

Niklas Christoffer Petersen ^{a,b,*}, Filipe Rodrigues ^b, Francisco Camara Pereira ^b

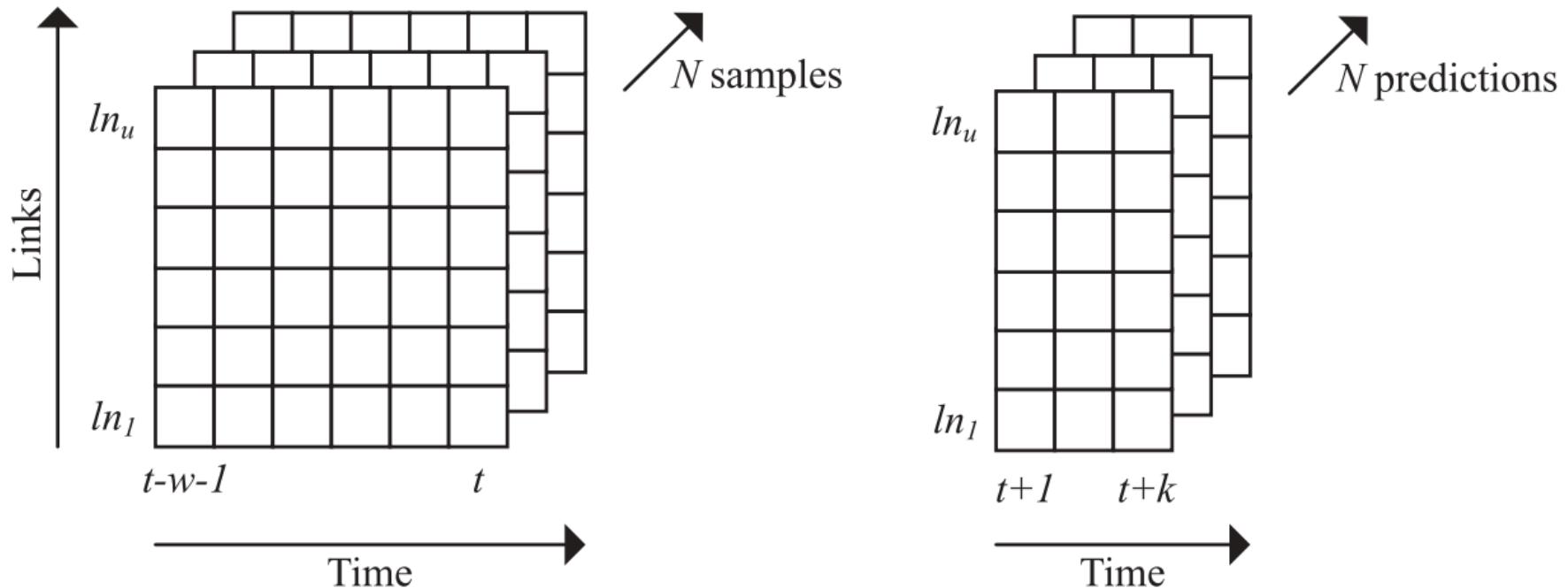


Bus travel time prediction

Table 1

Example of raw travel time measurements.

Timestamp	Linkref.	Link travel time (s)
2017-10-10 00:20:02	29848:1254	63
2017-10-10 00:21:07	1254:1255	65
2017-10-10 00:21:51	1255:10115	44
:	:	:

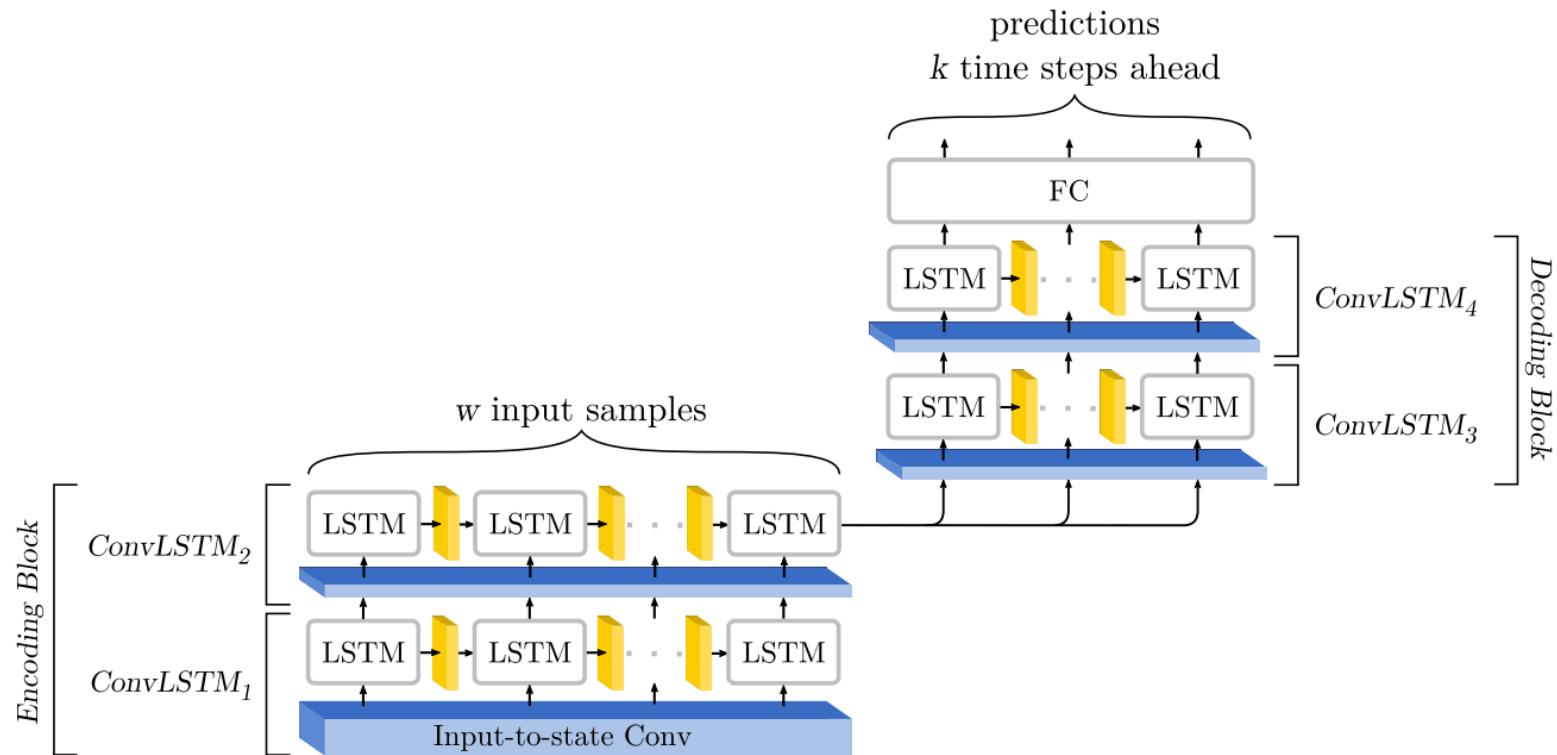


Bus travel time prediction

Table 1

Example of raw travel time measurements.

Timestamp	Linkref.	Link travel time (s)
2017-10-10 00:20:02	29848:1254	63
2017-10-10 00:21:07	1254:1255	65
2017-10-10 00:21:51	1255:10115	44
:	:	:



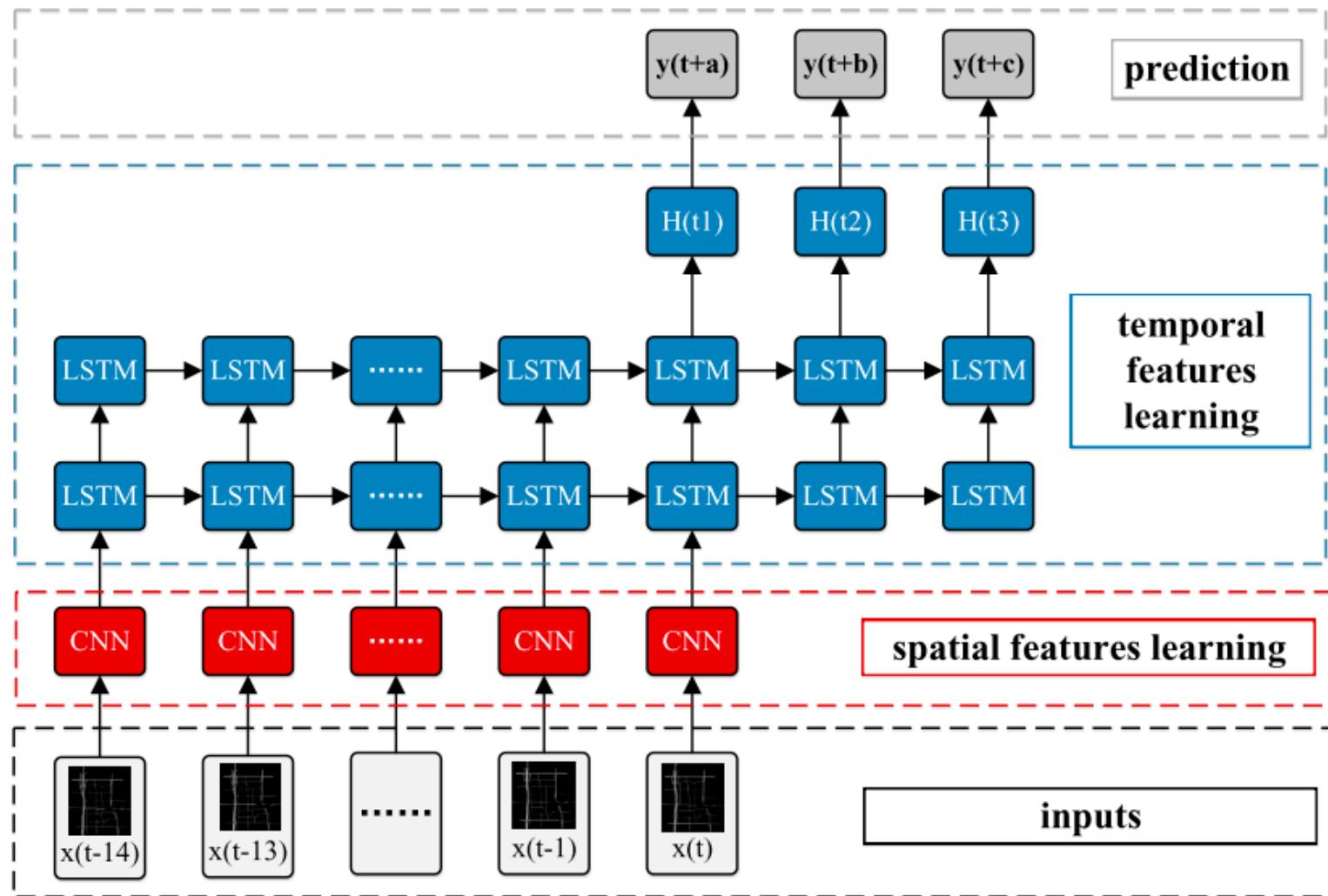
Bus travel time prediction

Table 2

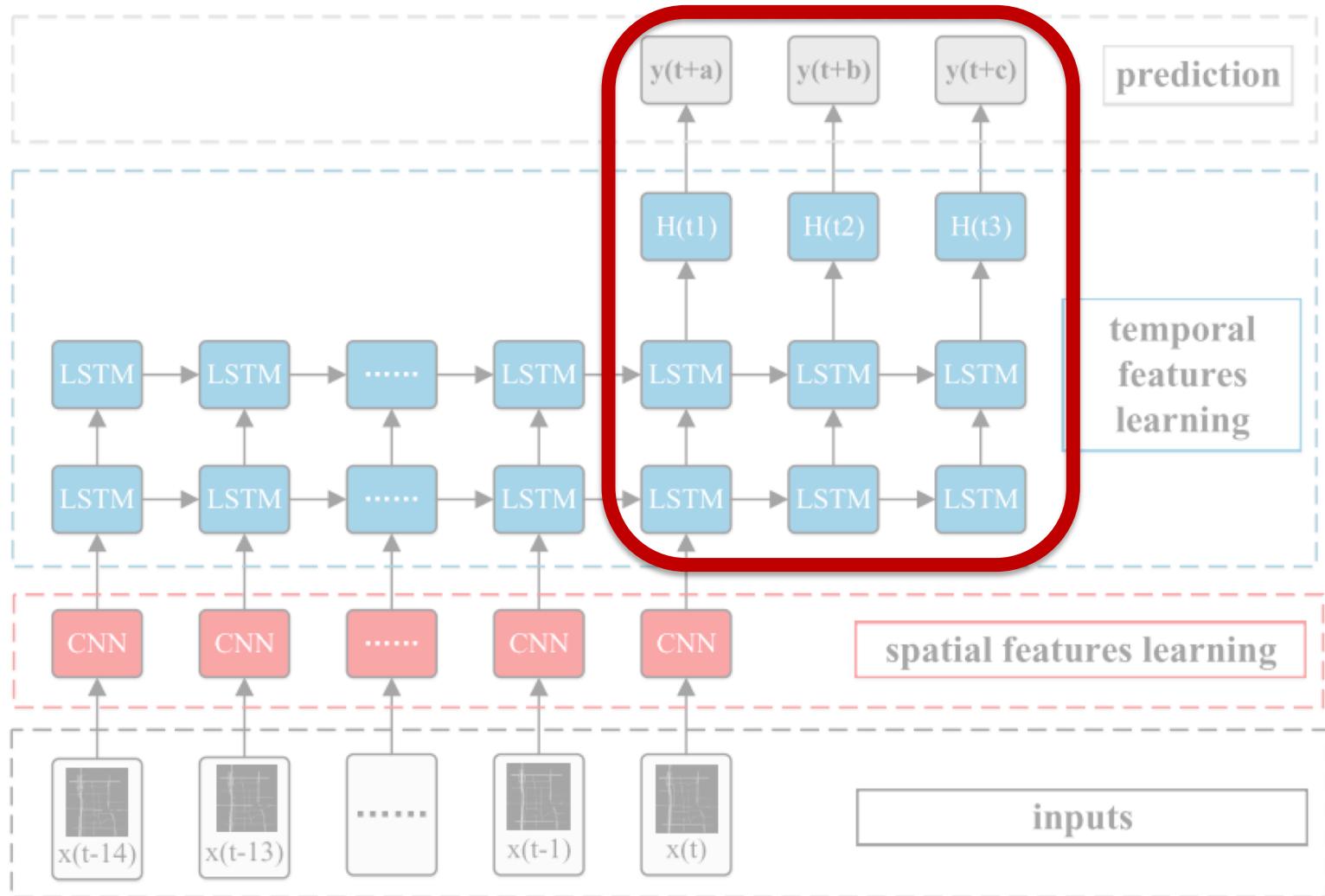
Results of the proposed and the baseline models.

Model	Time ahead	RMSE (min)	MAE (min)	MAPE (%)
Historical average		4.35	3.23	6.51%
Current model	$t + 1$ (15 min)	4.92	3.90	8.05%
	$t + 2$ (30 min)	4.91	3.46	6.82%
	$t + 3$ (45 min)	5.47	4.15	8.68%
Pure LSTM	$t + 1$ (15 min)	3.48	2.48	5.02%
	$t + 2$ (30 min)	3.56	2.51	5.08%
	$t + 3$ (45 min)	3.68	2.62	5.34%
Google Traffic	$t + 1$ (15 min)	3.67	2.96	6.32%
ConvLSTM	$t + 1$ (15 min)	2.66	1.99	4.19%
	$t + 2$ (30 min)	2.89	2.11	4.44%
	$t + 3$ (45 min)	3.11	2.27	4.75%

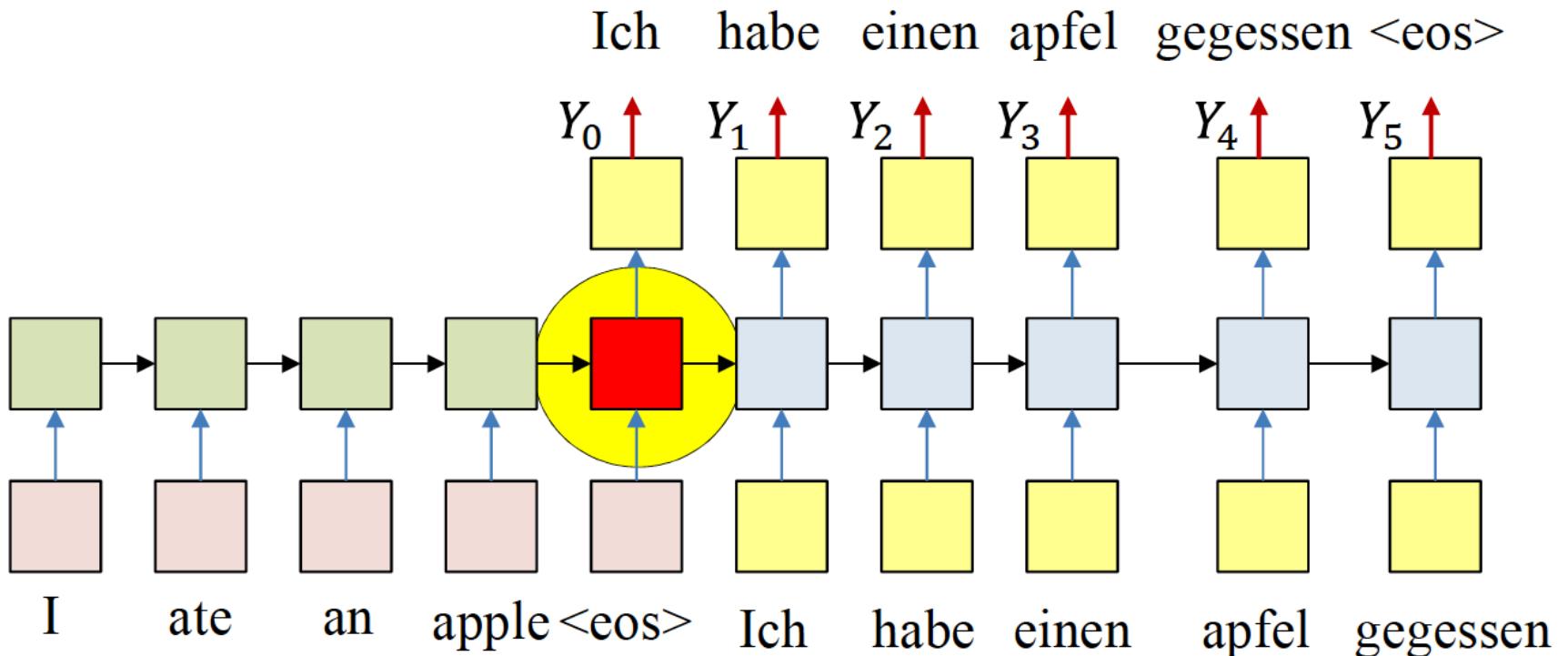
Two issues



Two issues

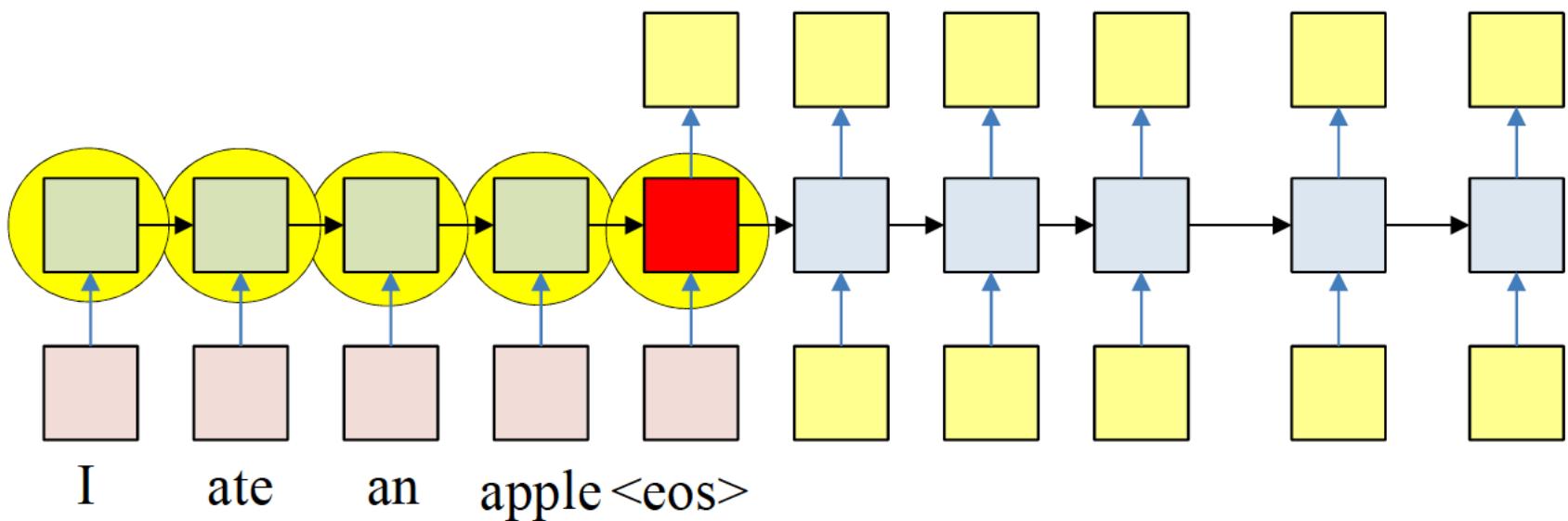


Neural machine translation



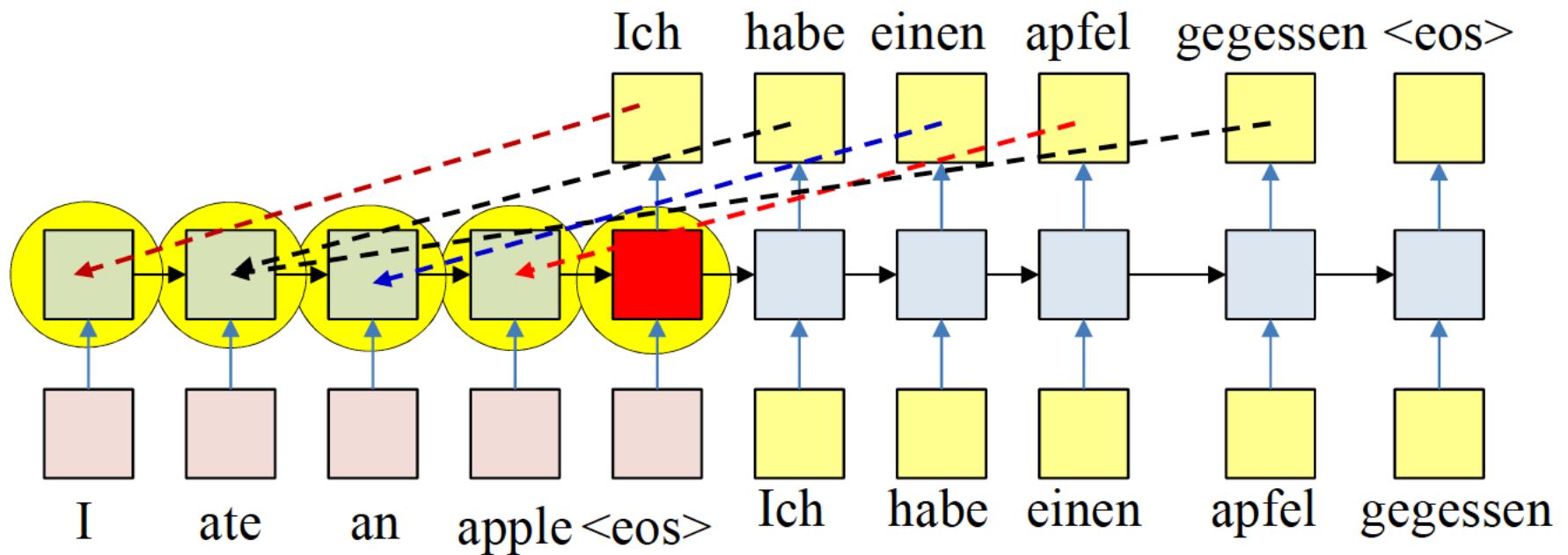
Source: Bhiksha Raj

Neural machine translation



Source: Bhiksha Raj

Neural machine translation



Source: Bhiksha Raj

2014

(dramatic reenactment)

2014

*Neural
Machine
Translation*

MT research

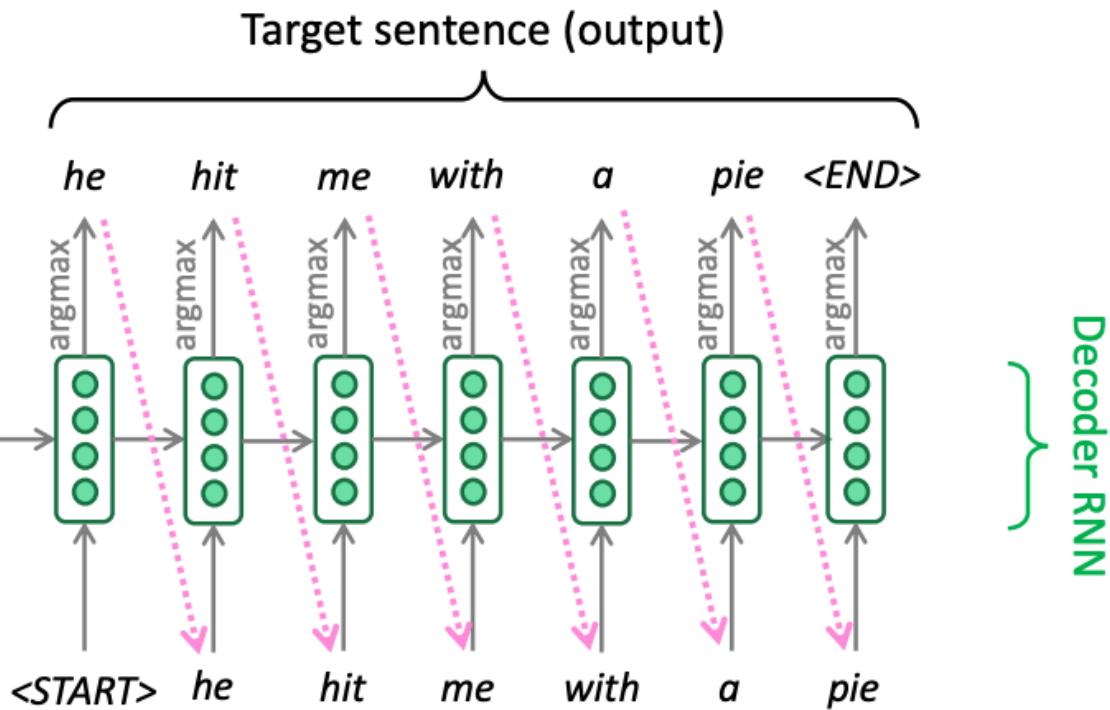
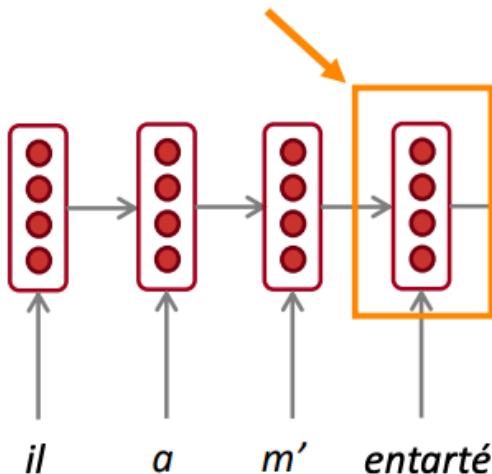
(dramatic reenactment)

Neural machine translation

The sequence-to-sequence model

Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Encoder RNN



Source: Abigail See

Attention

- Attend to different parts of the input based on the output

$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'=1}^S \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad [\text{Attention weights}] \quad (1)$$

$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s \quad [\text{Context vector}] \quad (2)$$

$$\mathbf{a}_t = f(\mathbf{c}_t, \mathbf{h}_t) = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad [\text{Attention vector}] \quad (3)$$

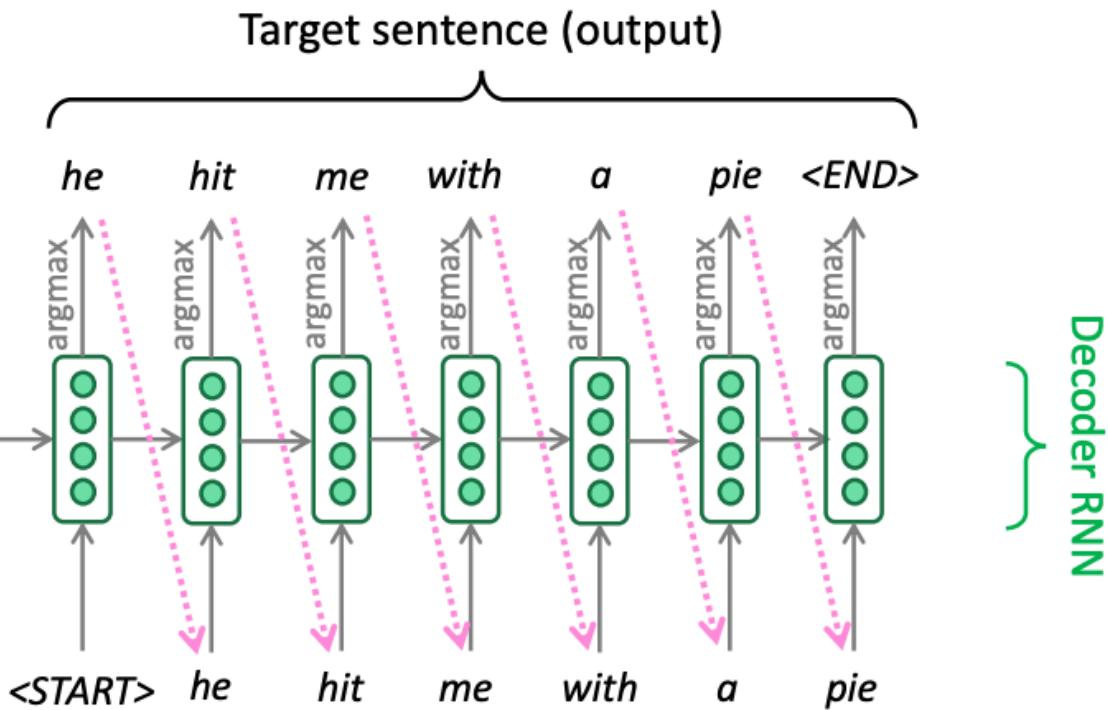
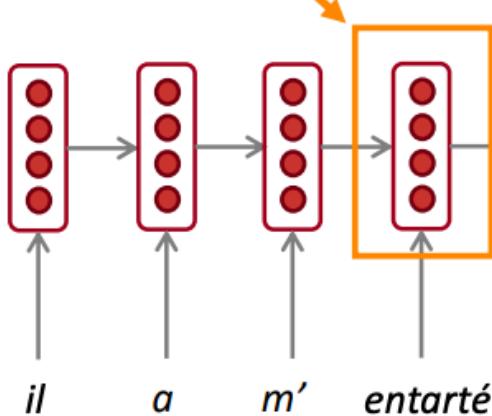
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \mathbf{W} \bar{\mathbf{h}}_s & [\text{Luong's multiplicative style}] \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_1 \mathbf{h}_t + \mathbf{W}_2 \bar{\mathbf{h}}_s) & [\text{Bahdanau's additive style}] \end{cases} \quad (4)$$

Neural machine translation

The sequence-to-sequence model

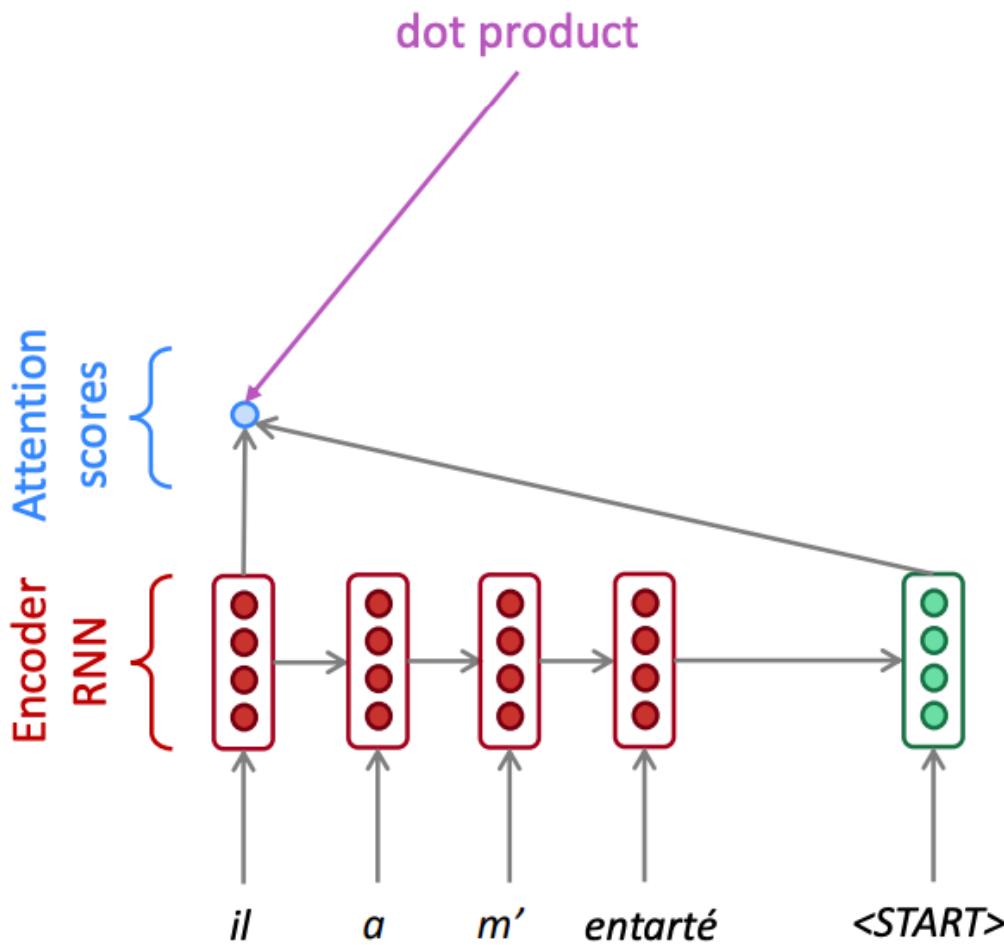
Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Encoder RNN



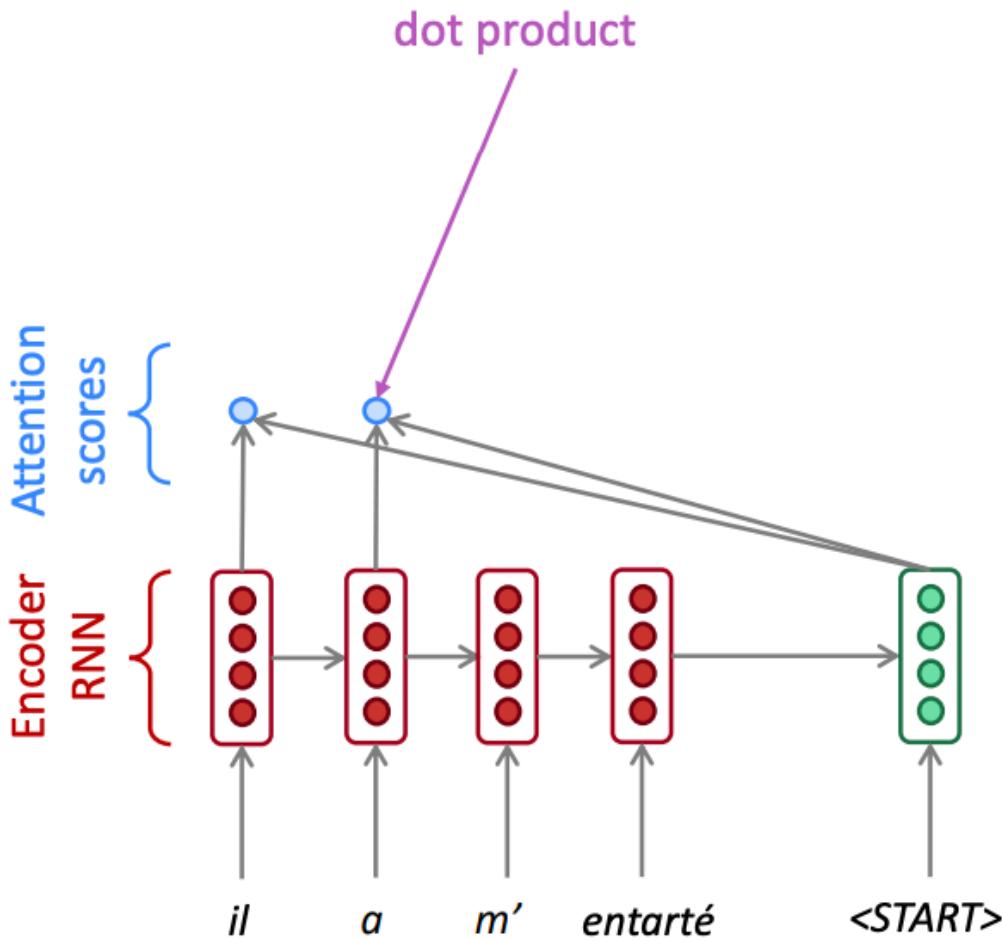
Source: Abigail See

Neural machine translation



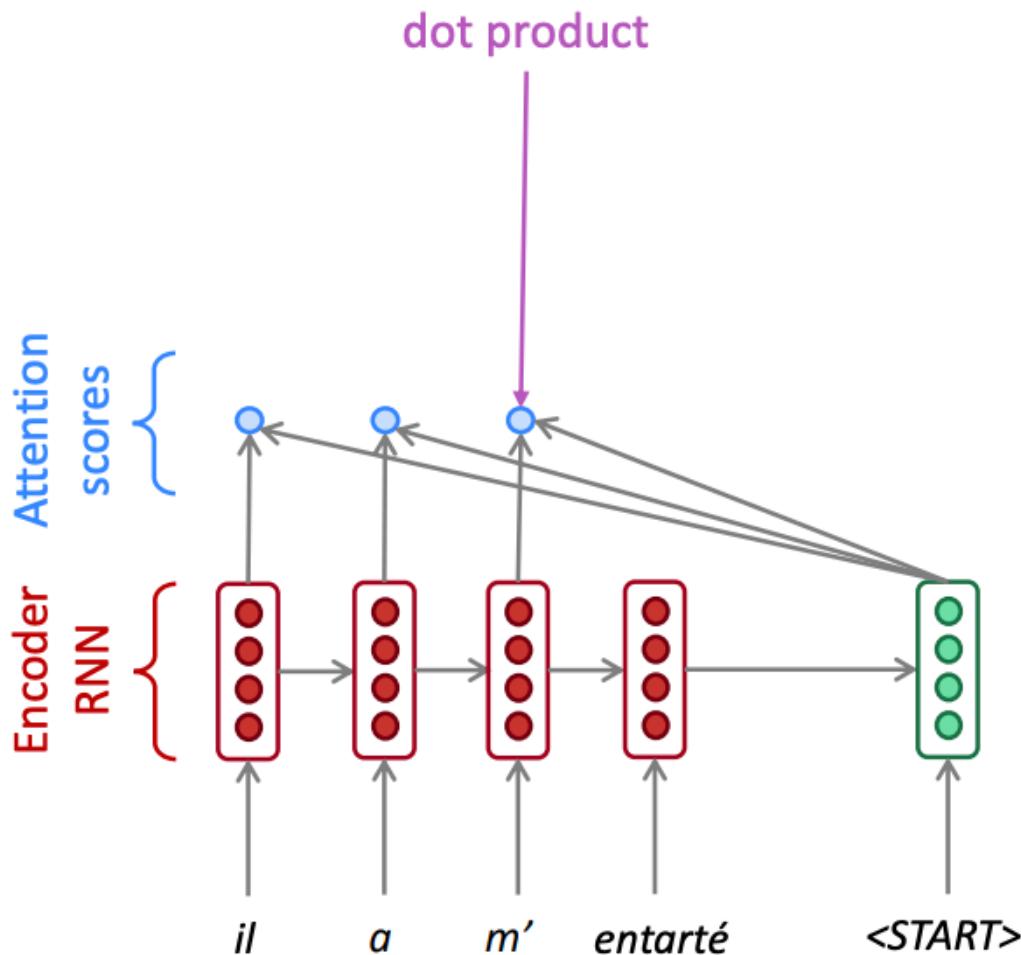
Source: Abigail See

Neural machine translation



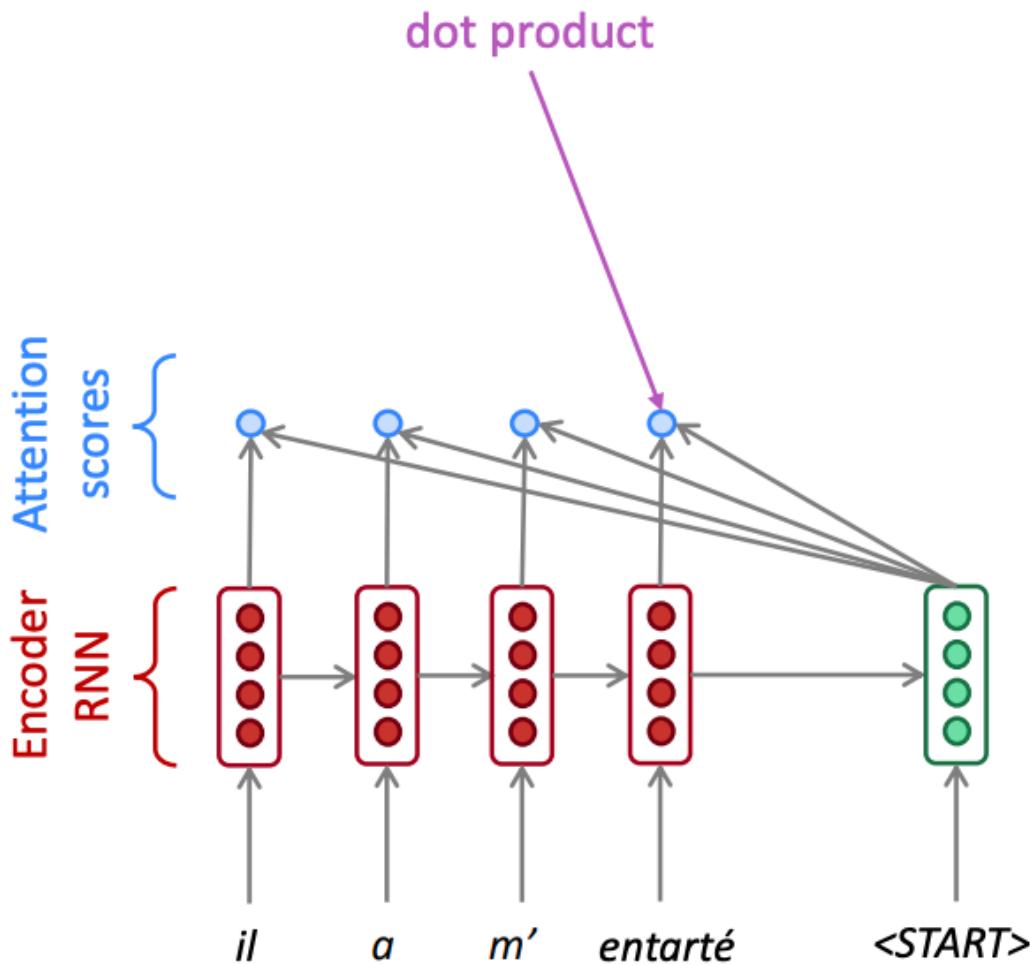
Source: Abigail See

Neural machine translation



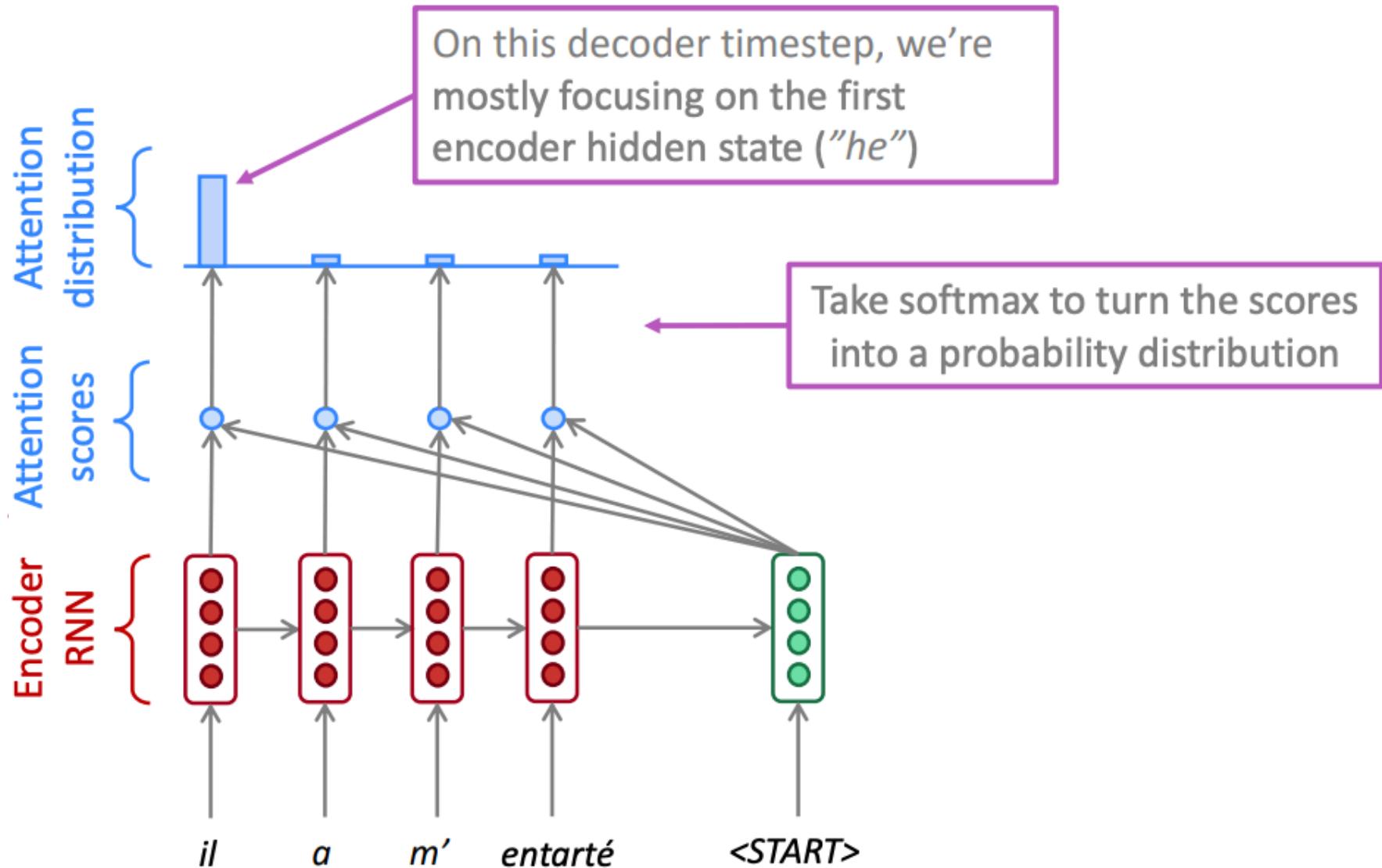
Source: Abigail See

Neural machine translation

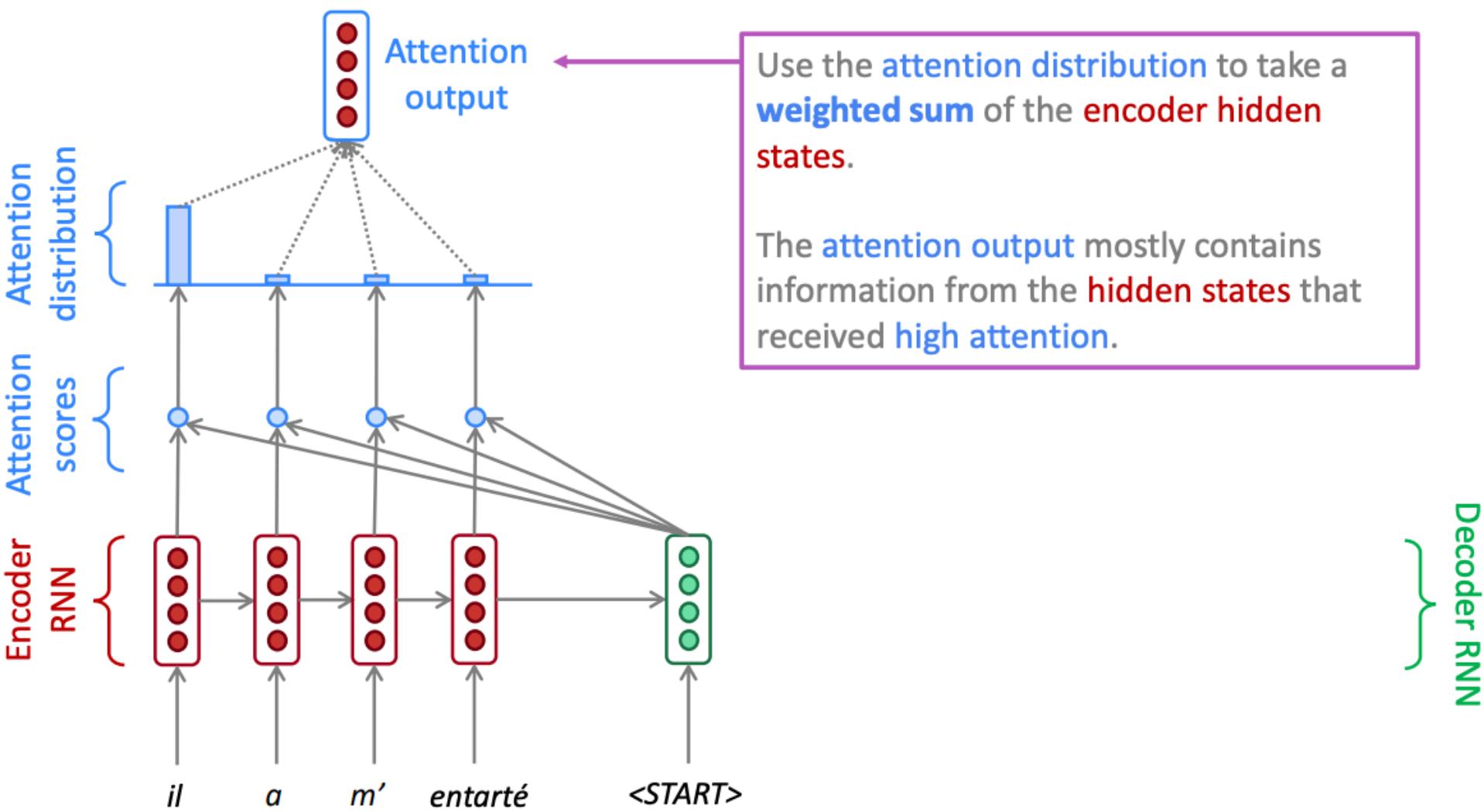


Source: Abigail See

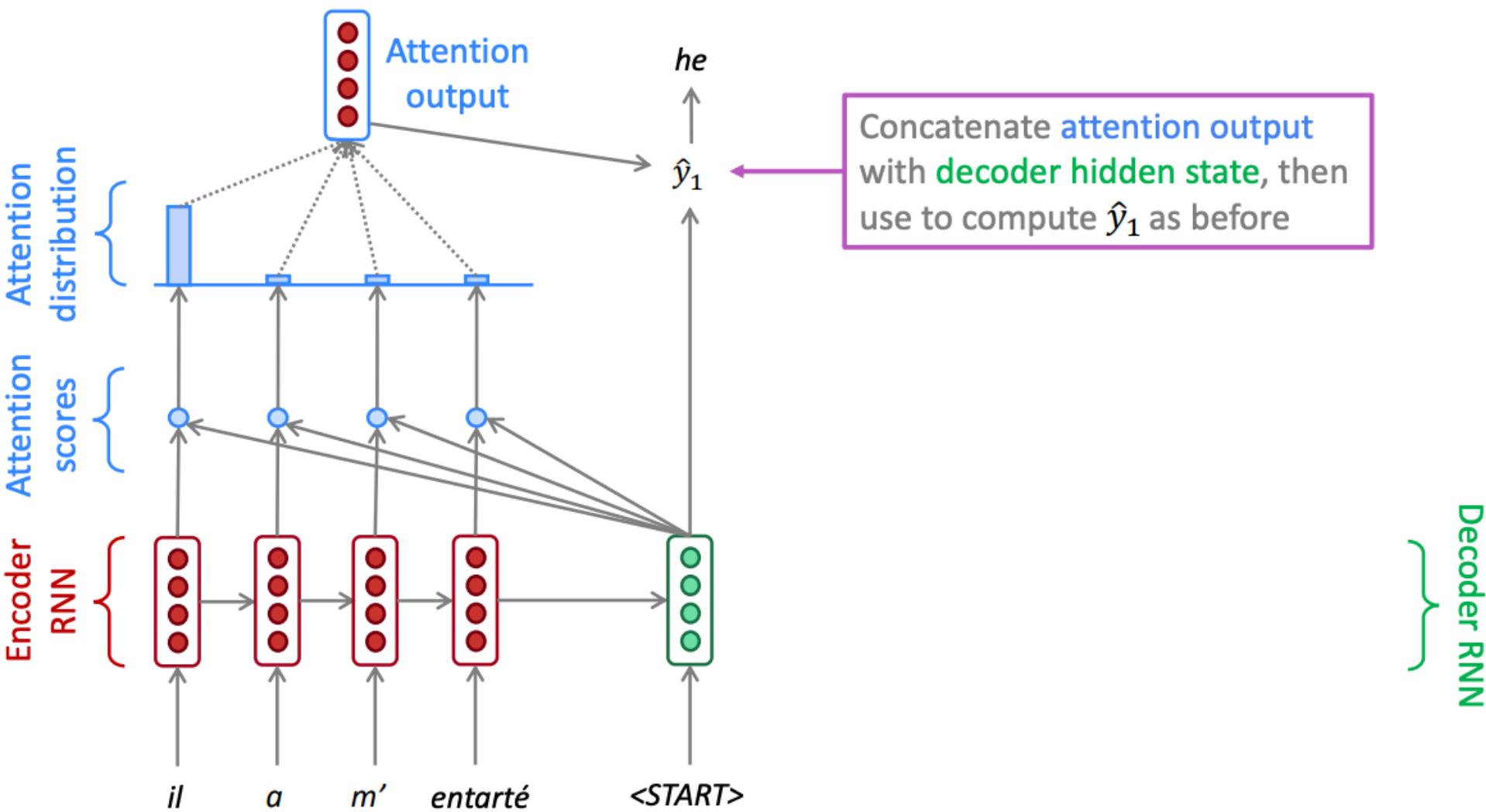
Neural machine translation



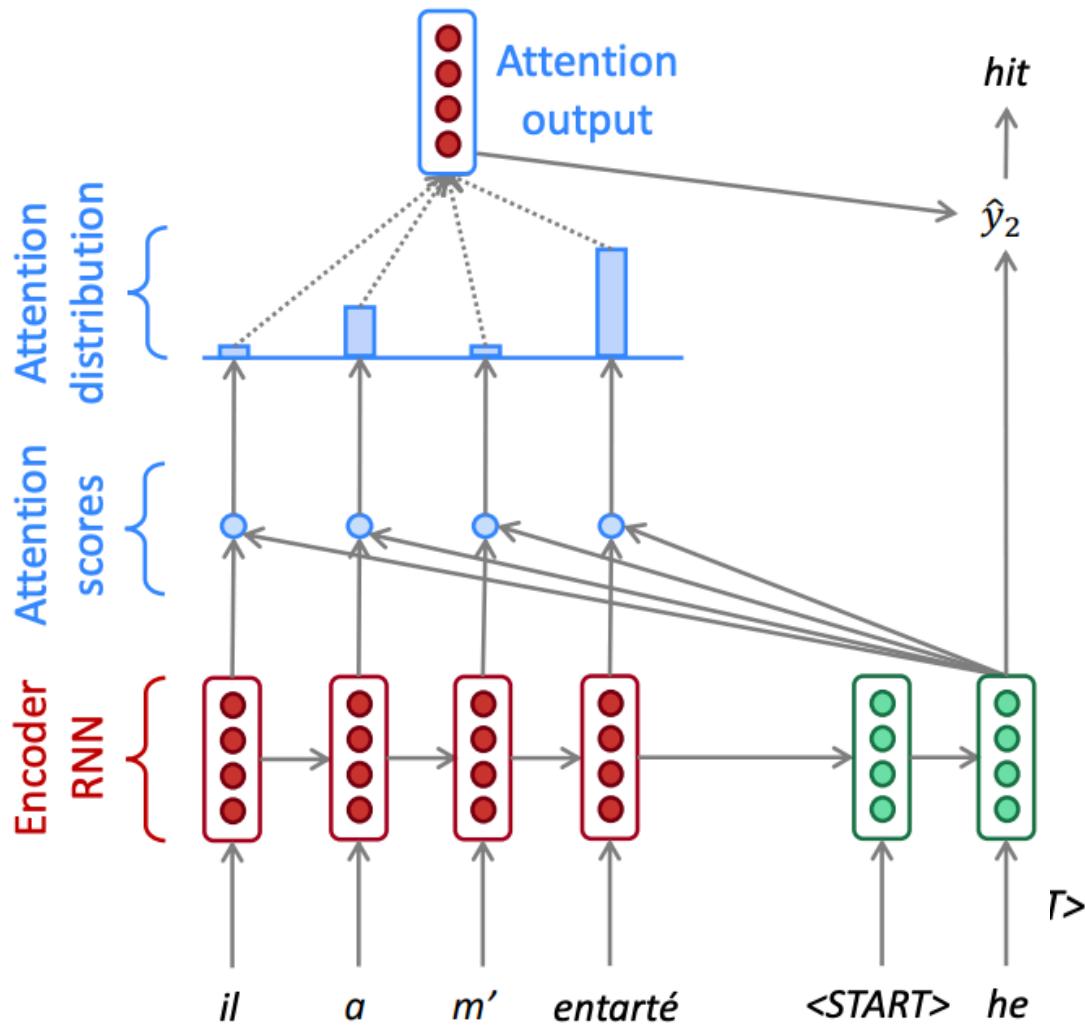
Neural machine translation



Neural machine translation

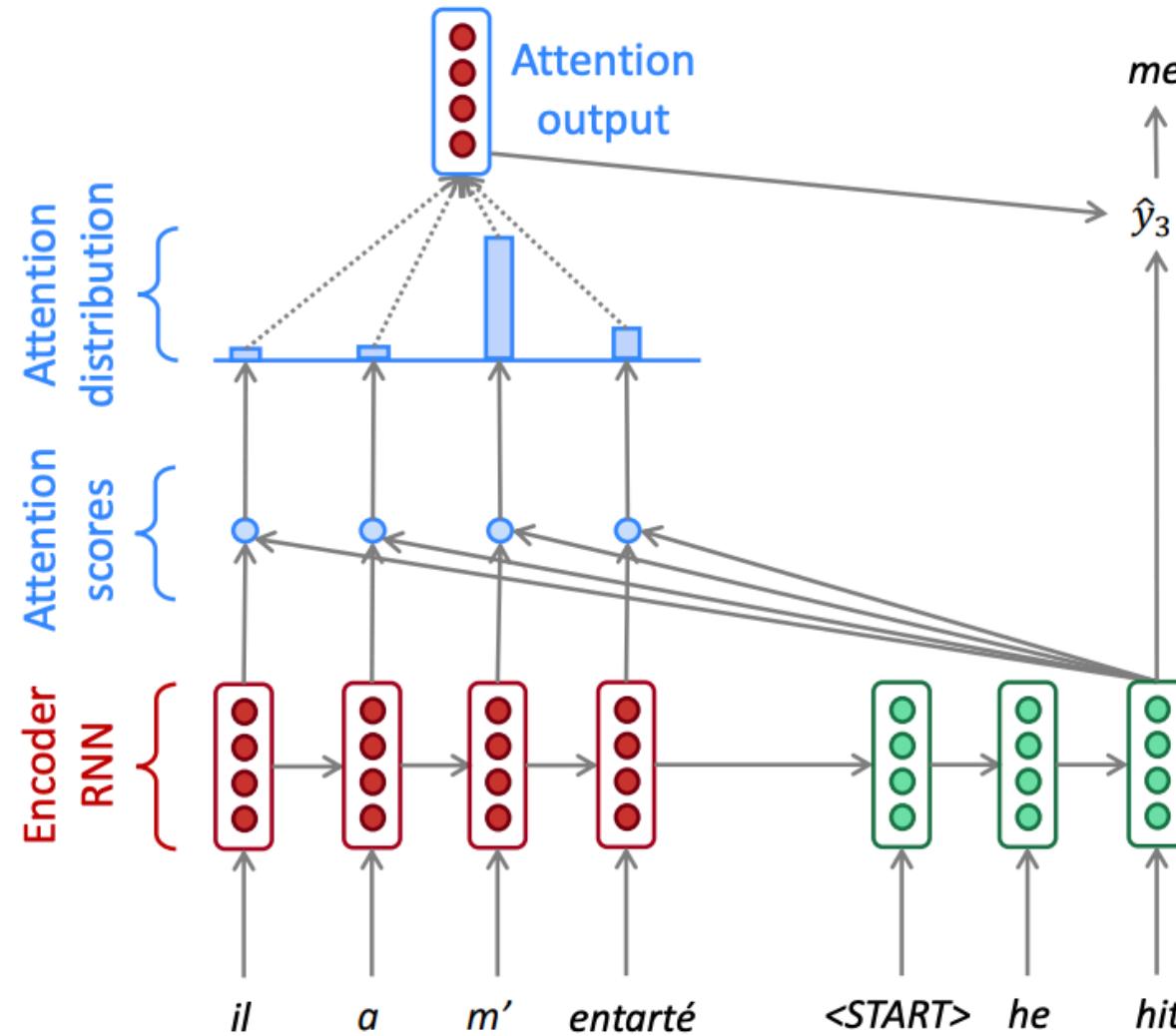


Neural machine translation



Source: Abigail See

Neural machine translation



Source: Abigail See

Lab 2

Keras

TensorFlow / Theano / CNTK / ...

CUDA / cuDNN

BLAS, Eigen

GPU

CPU

Individual mobility prediction

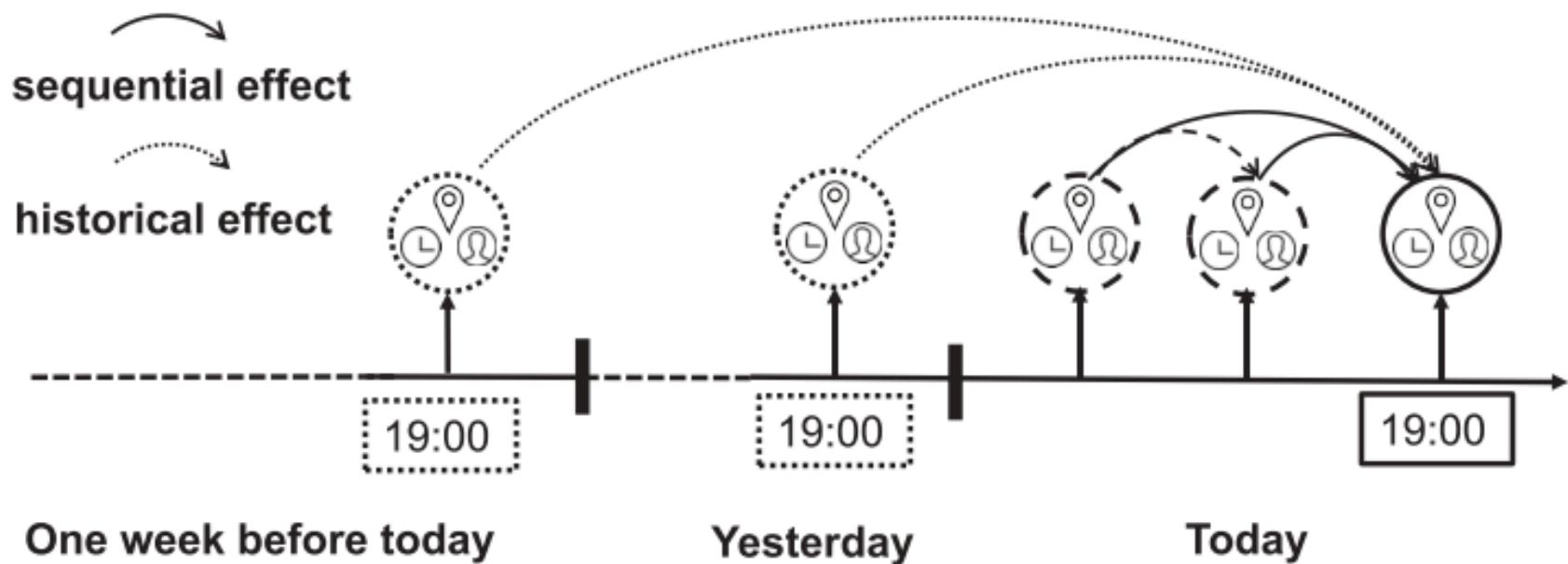
DEFINITION 1 (TRAJECTORY SEQUENCE). *Spatiotemporal point q is a tuple of time stamp t and location identification l , i.e., $q = (t, l)$. Given a user identification u , trajectory sequence S is a spatiotemporal point sequence, i.e., $S^u = q_1 q_2 \dots q_n$.*

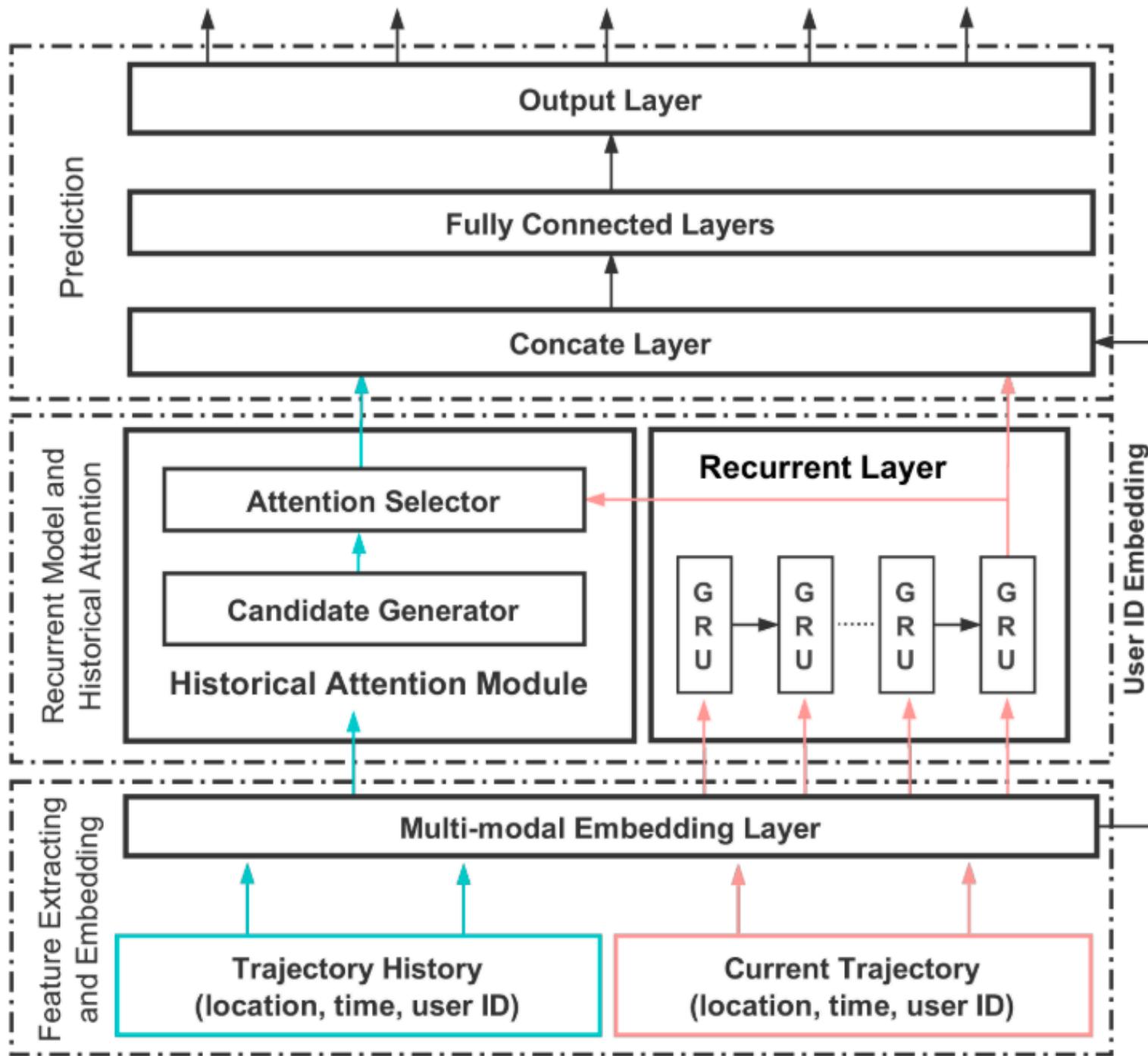
DEFINITION 2 (TRAJECTORY). *Given a trajectory sequence S^u and a time window t_w , trajectory is a subsequence $S_{t_w}^u = q_i q_{i+1} \dots q_{i+k}$ of S^u in the time window t_w , if $\forall 1 < j \leq k$, t_{q_j} belongs to t_w .*

DeepMove: Predicting Human Mobility with Attentional Recurrent Networks

Individual mobility prediction

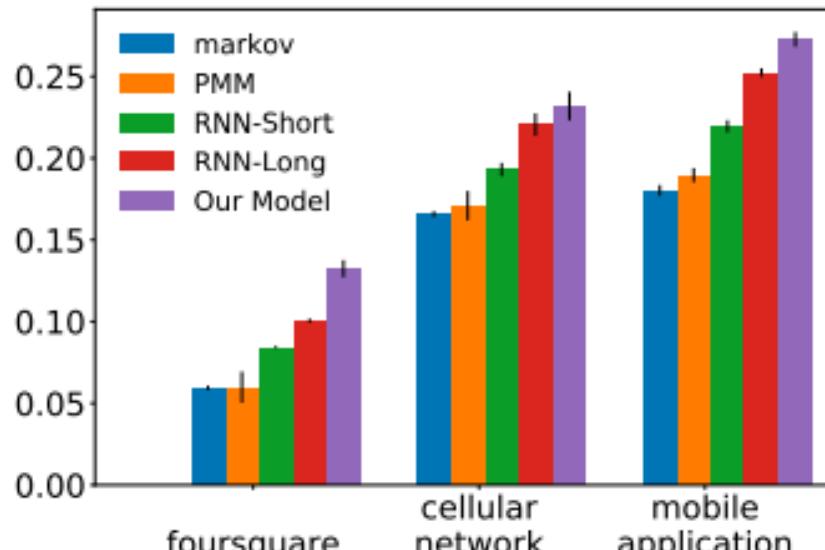
PROBLEM 1 (MOBILITY PREDICTION). *Given the current trajectory $S_{t_{w_m}}^u = q_1 q_2 \dots q_n$ and the corresponding trajectory history $S_{t_{w_1}}^u S_{t_{w_2}}^u \dots S_{t_{w_{m-1}}}^u$, predict the next spatiotemporal point q_{n+1} in the trajectory.*



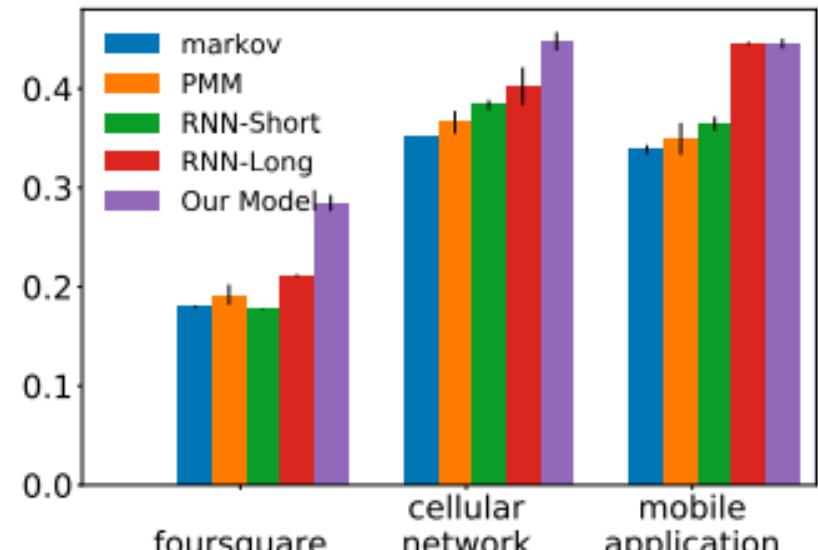


Individual mobility prediction

- Rank the candidate locations by the probabilities generated from the model, and check whether the ground-truth location v appears in the top- k candidate locations.

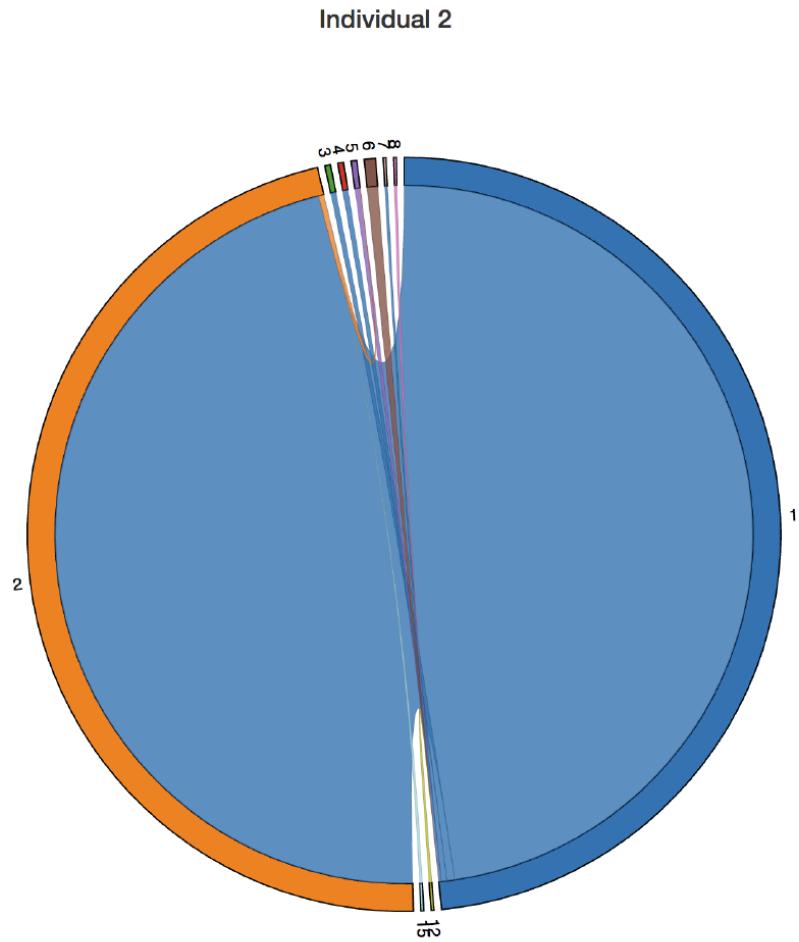
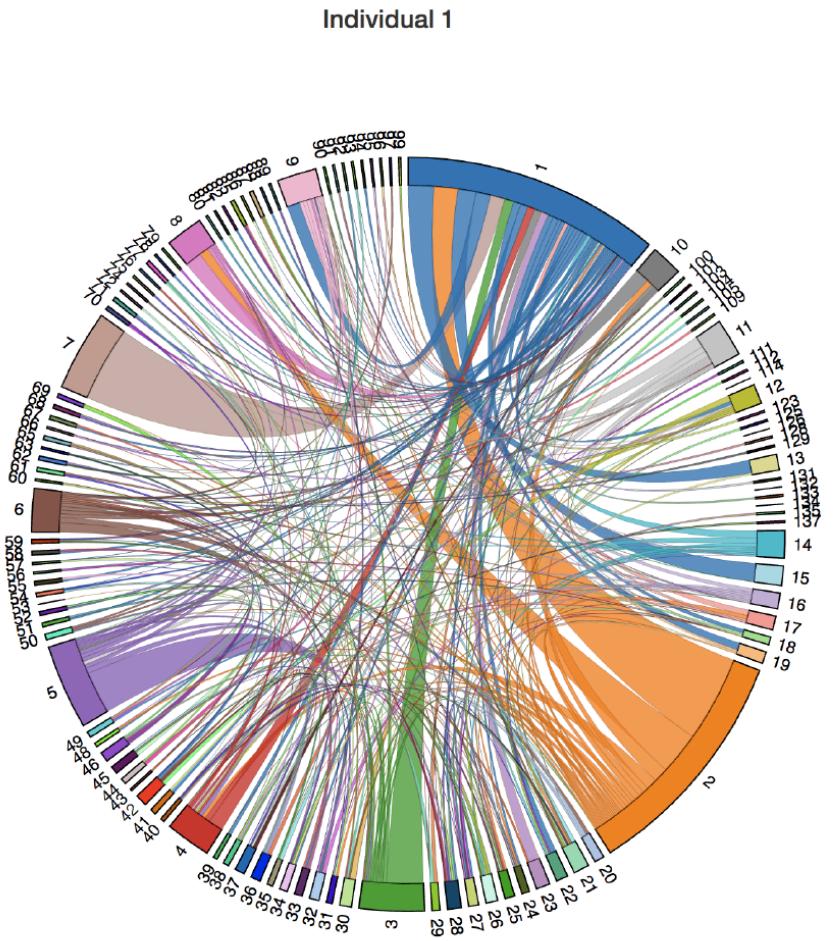


(a) top-1 prediction accuracy



(b) top-5 prediction accuracy

Individual mobility prediction



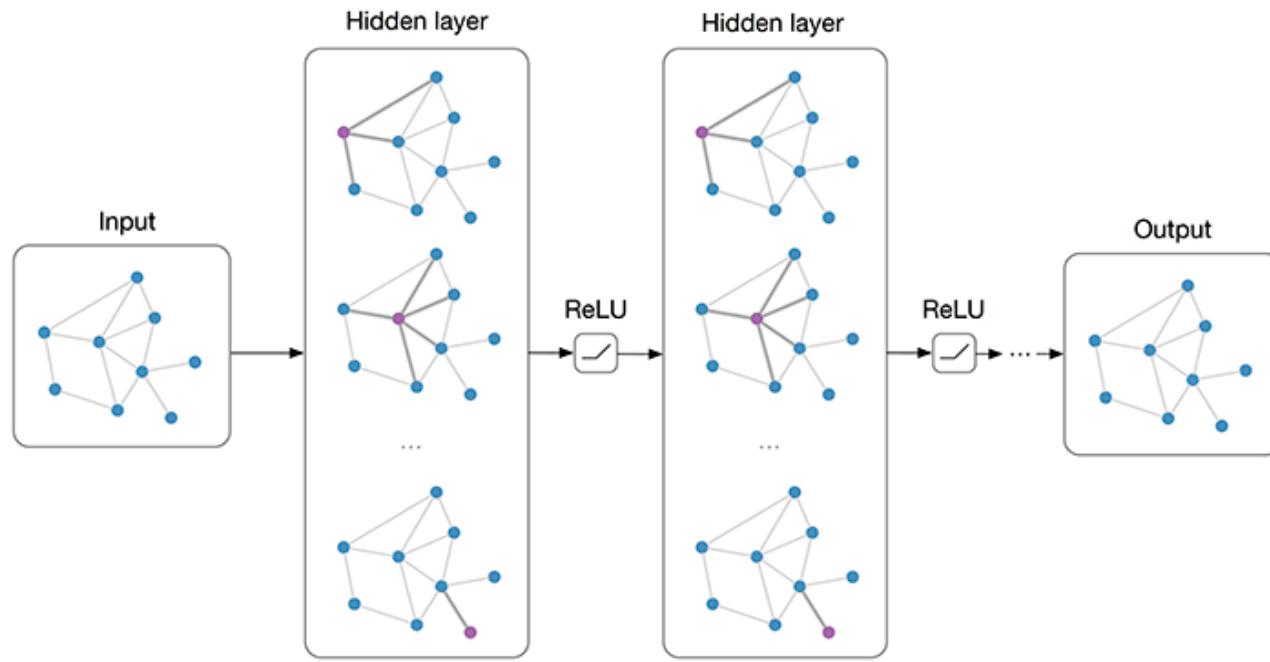
Graph Neural Networks

Graph Neural Networks

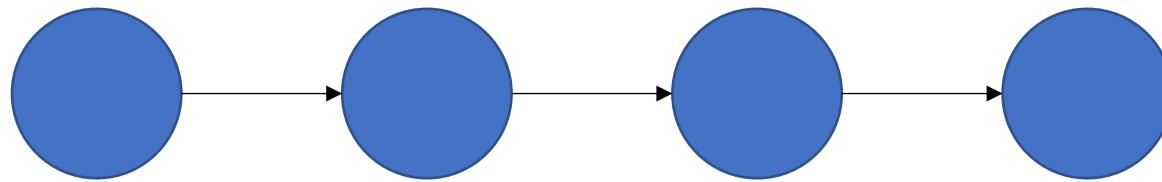
- CNN only works for regular grids

Graph Neural Networks

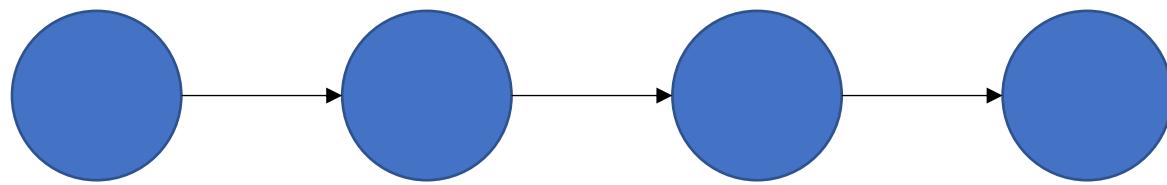
- CNN only works for regular grids
- We want something similar, but for irregular grids
 - Capture local dependencies
 - Translation invariant
 - Parameter sharing



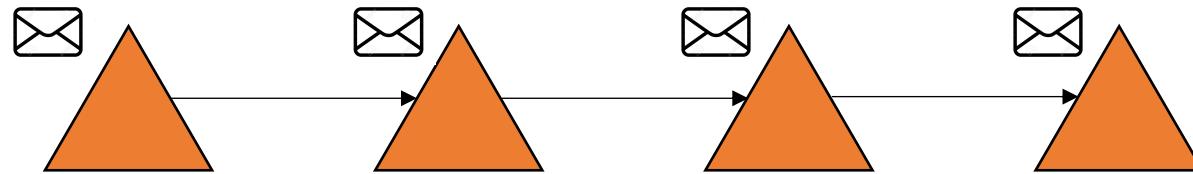
Neural message passing networks

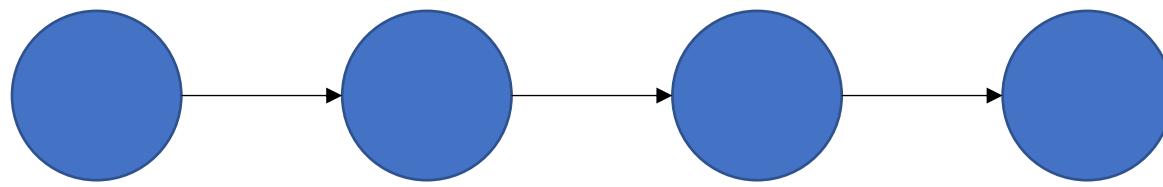


Embed('word') = 

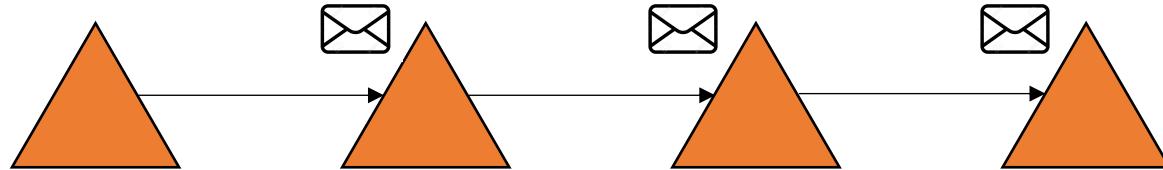


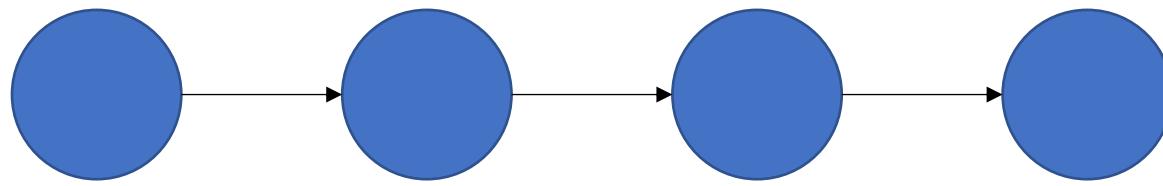
Embed('word') =



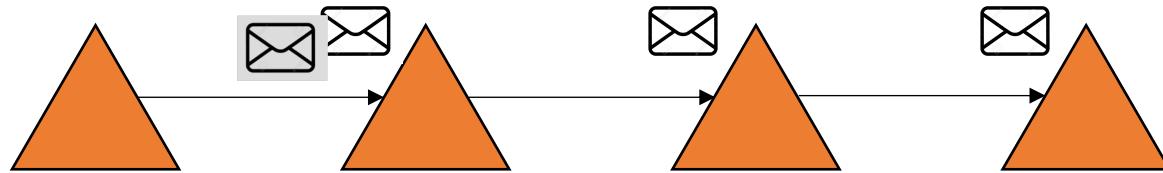


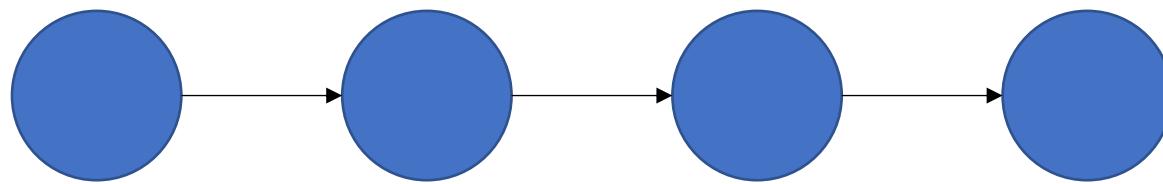
Embed('word') =



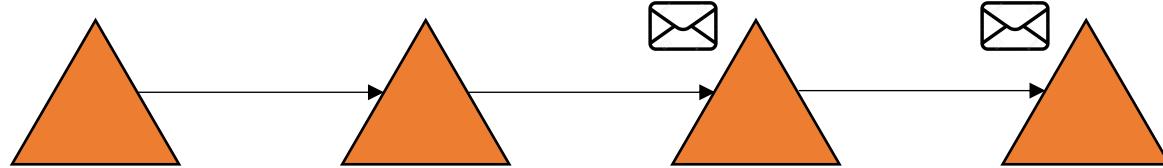


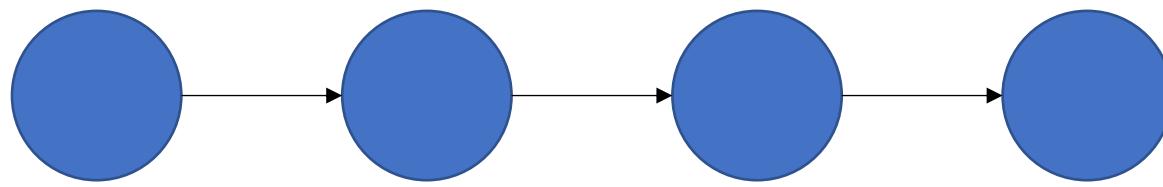
Embed('word') = 



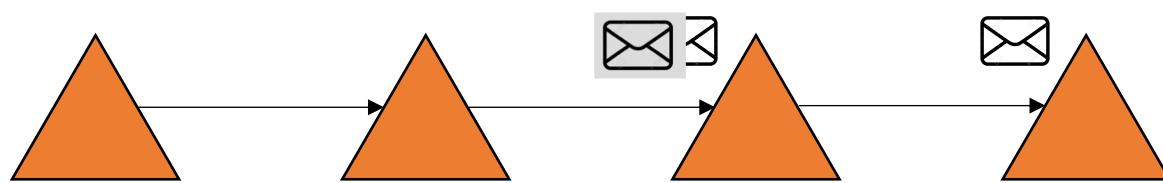


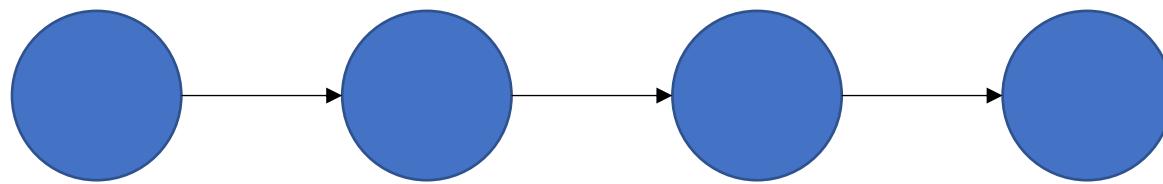
Embed('word') =





Embed('word') =

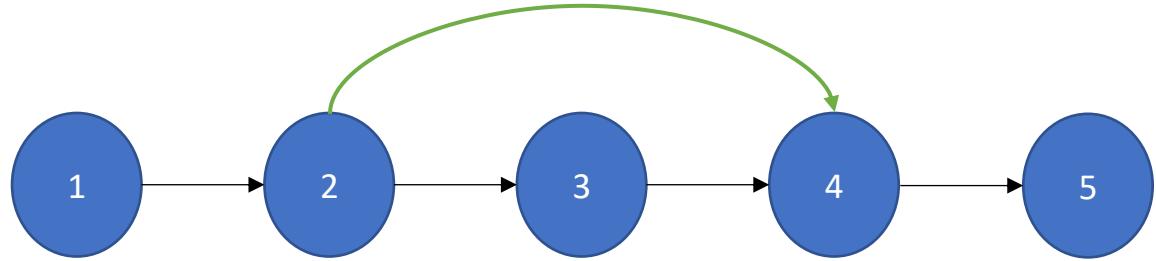
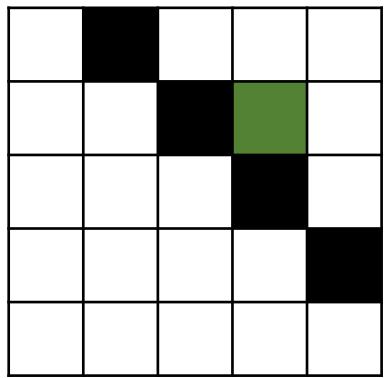


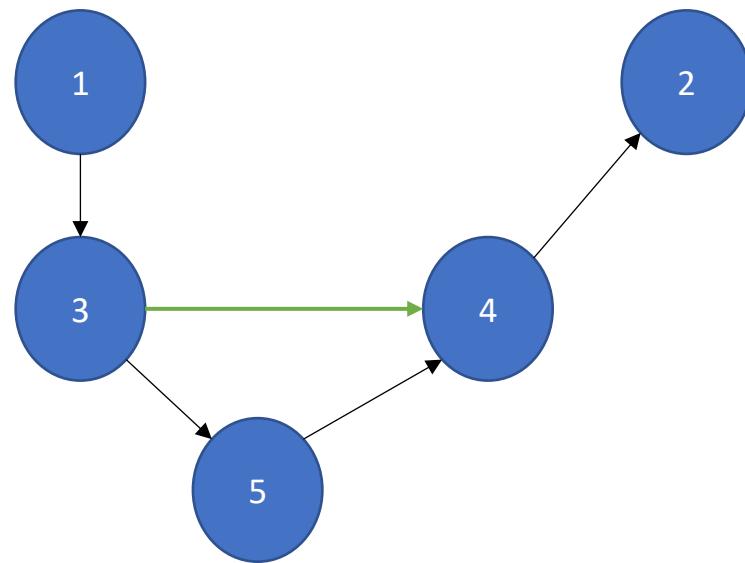
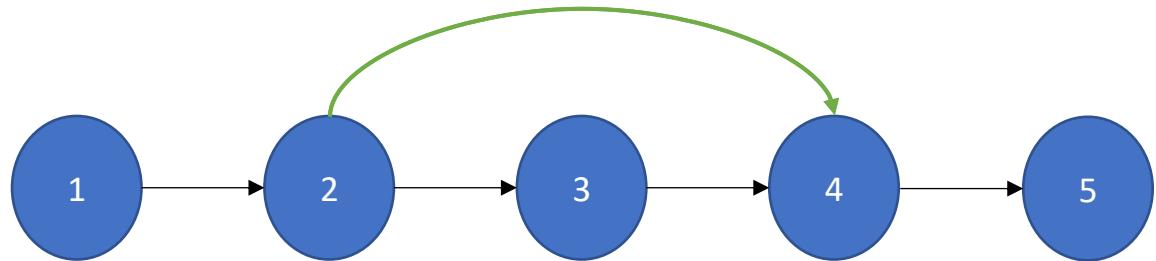
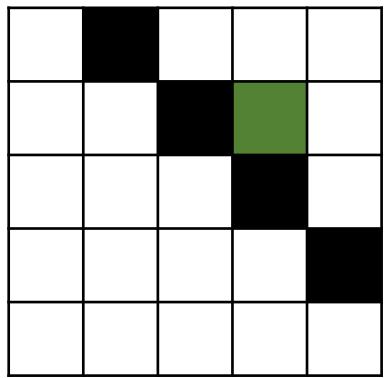


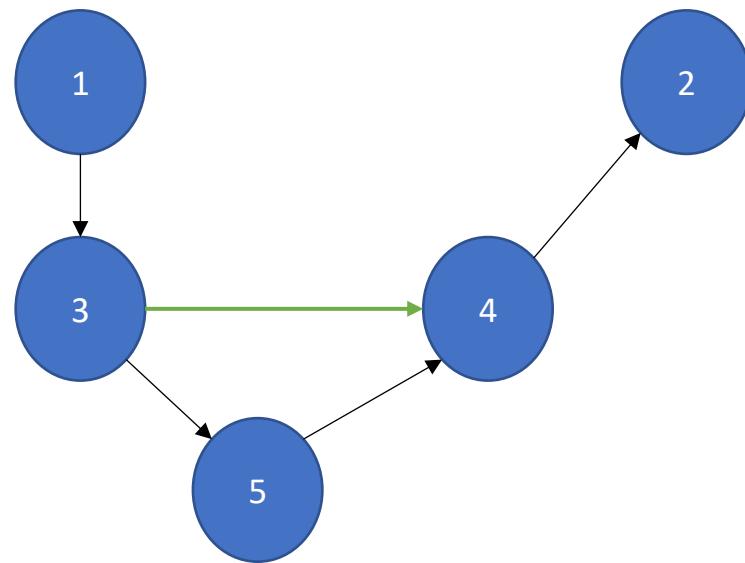
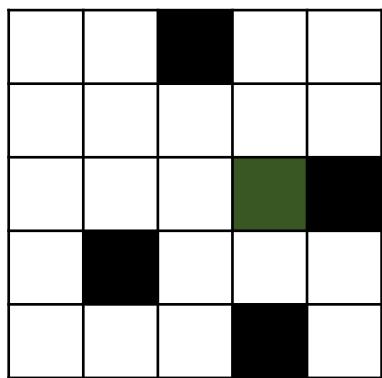
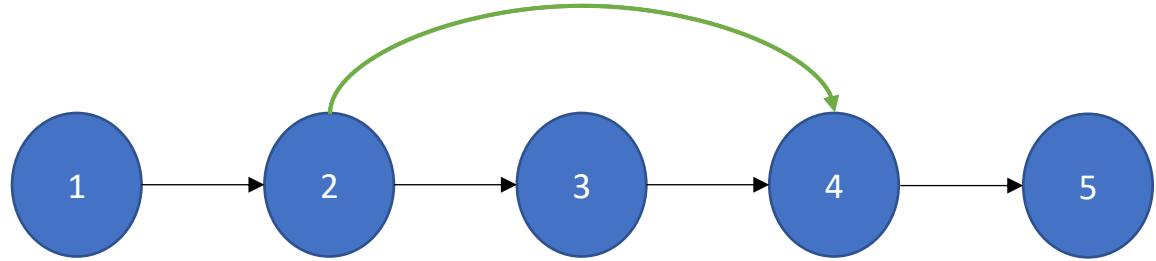
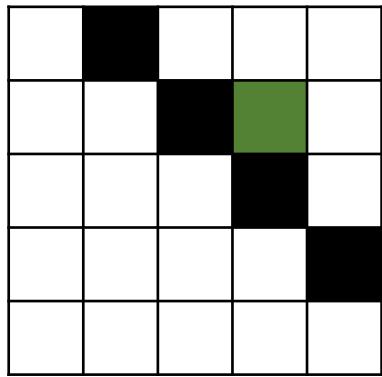
Embed('word') = 



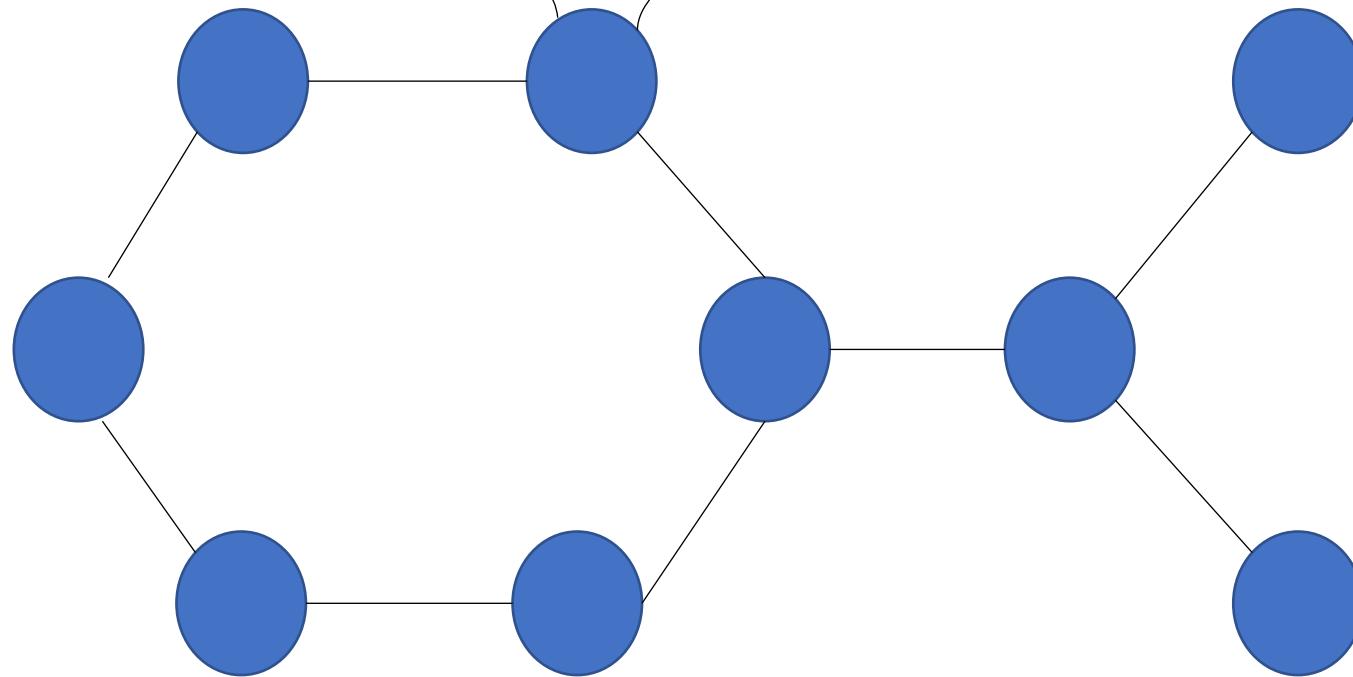
$$\text{✉} = \triangle (\text{✉}, \text{✉})$$

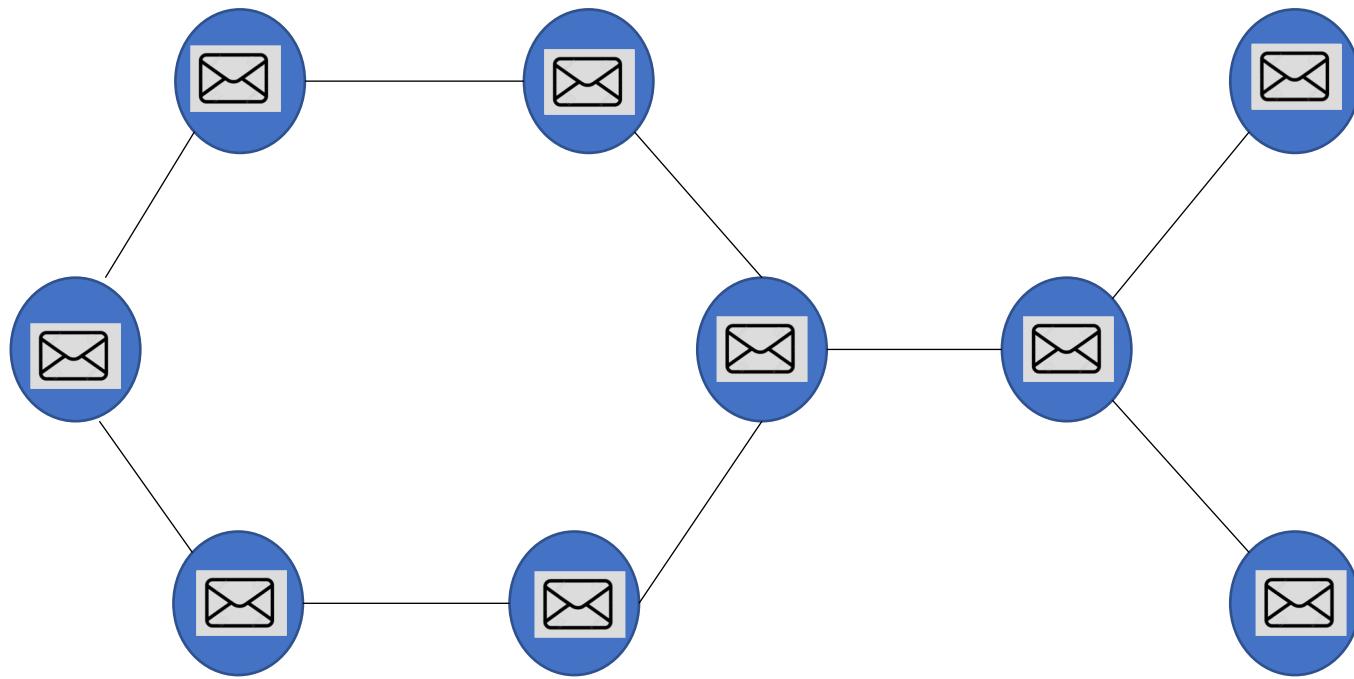


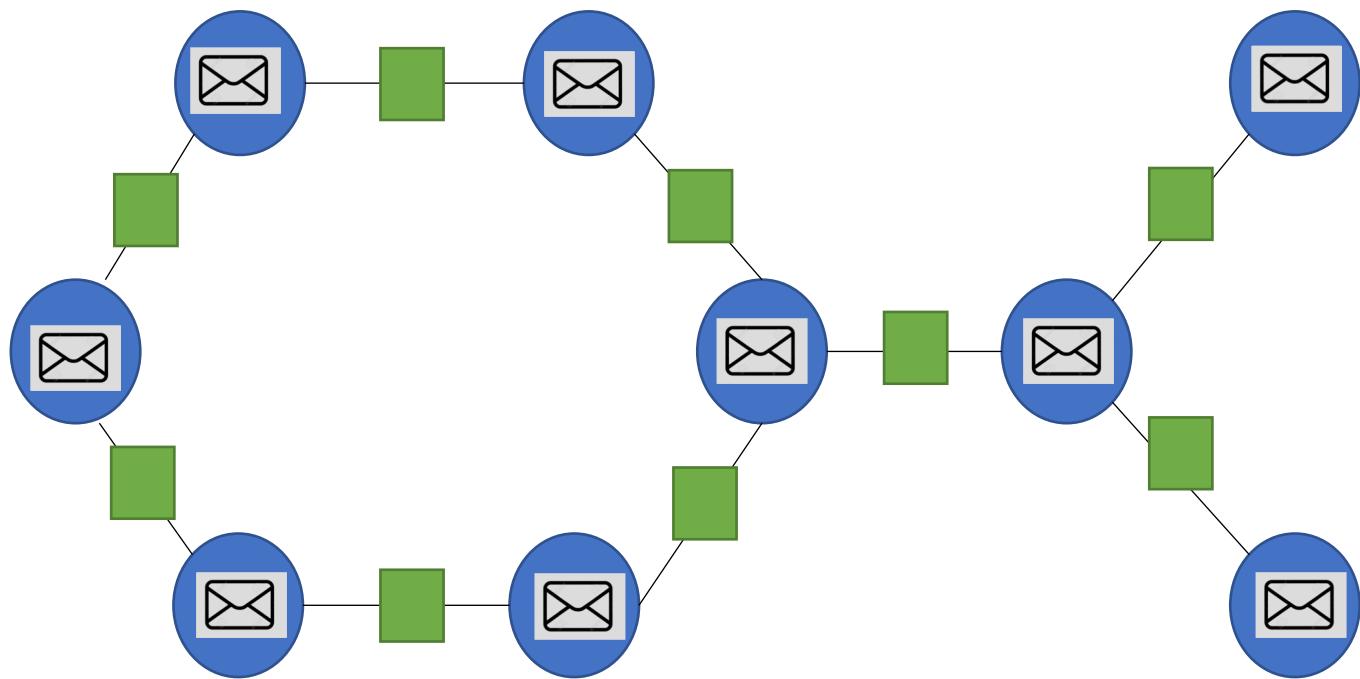


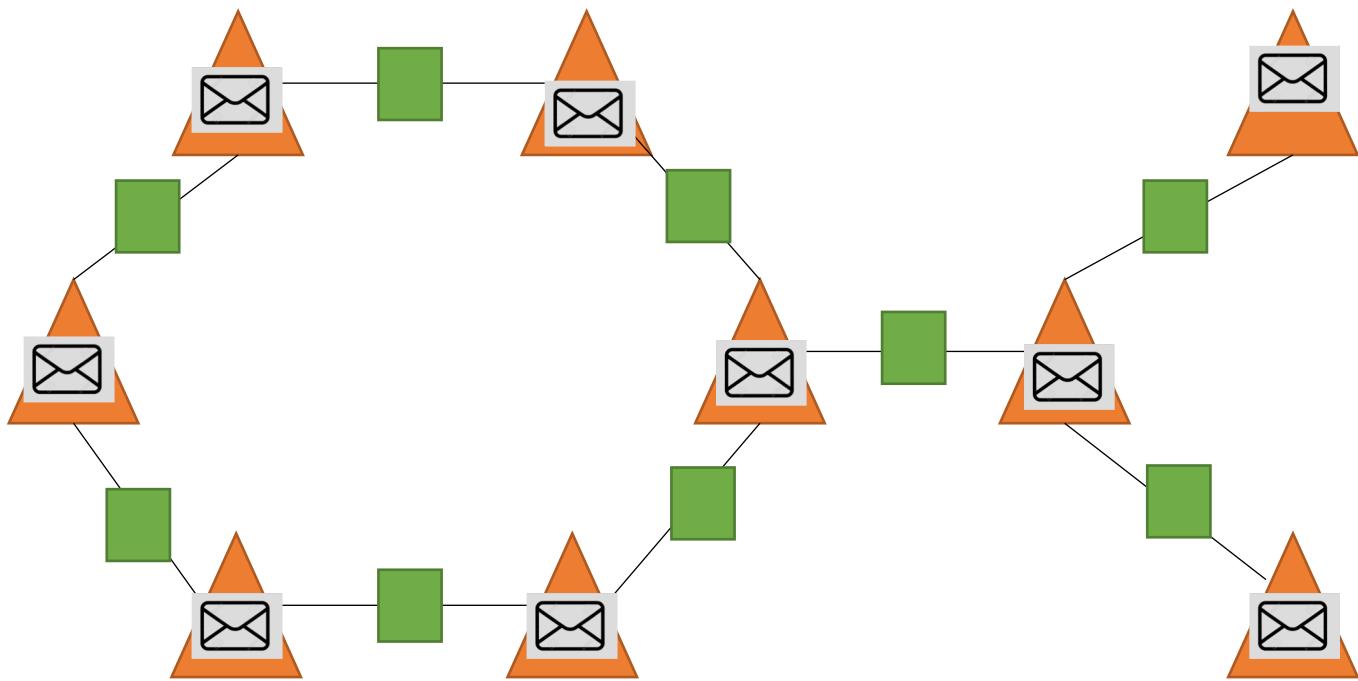


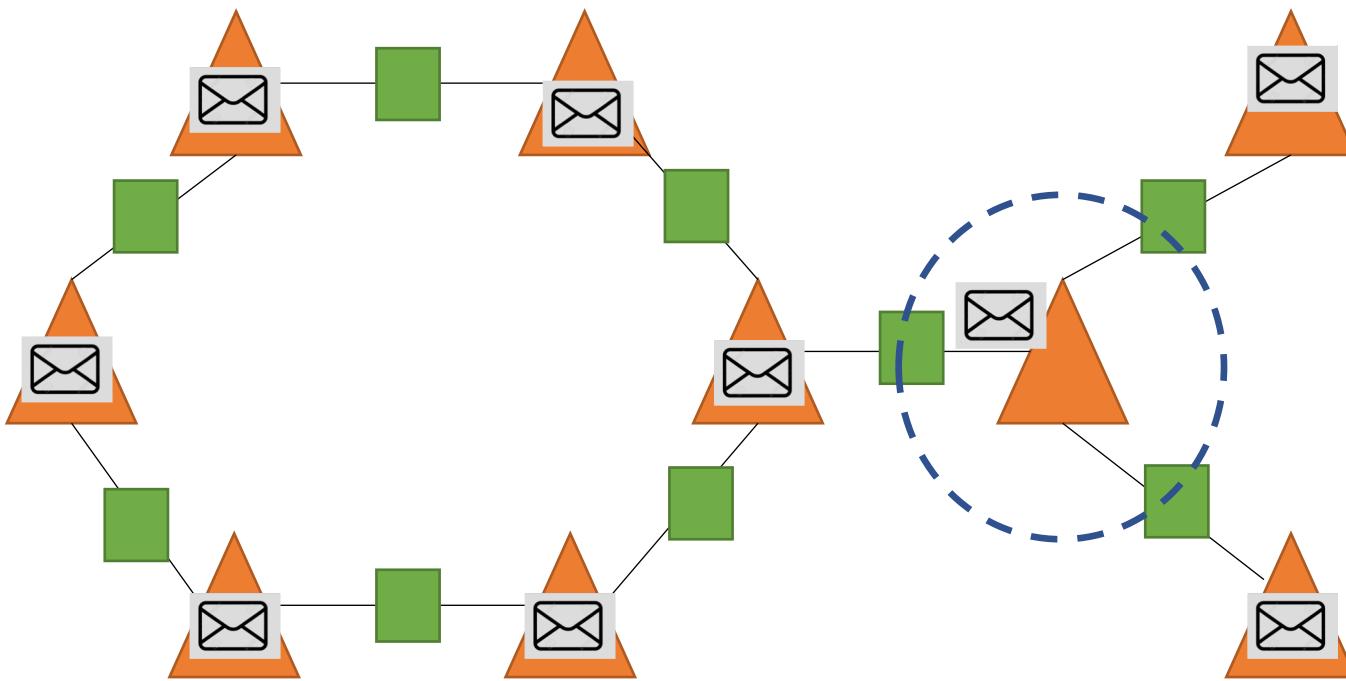
[25, 123, 148, 254] =

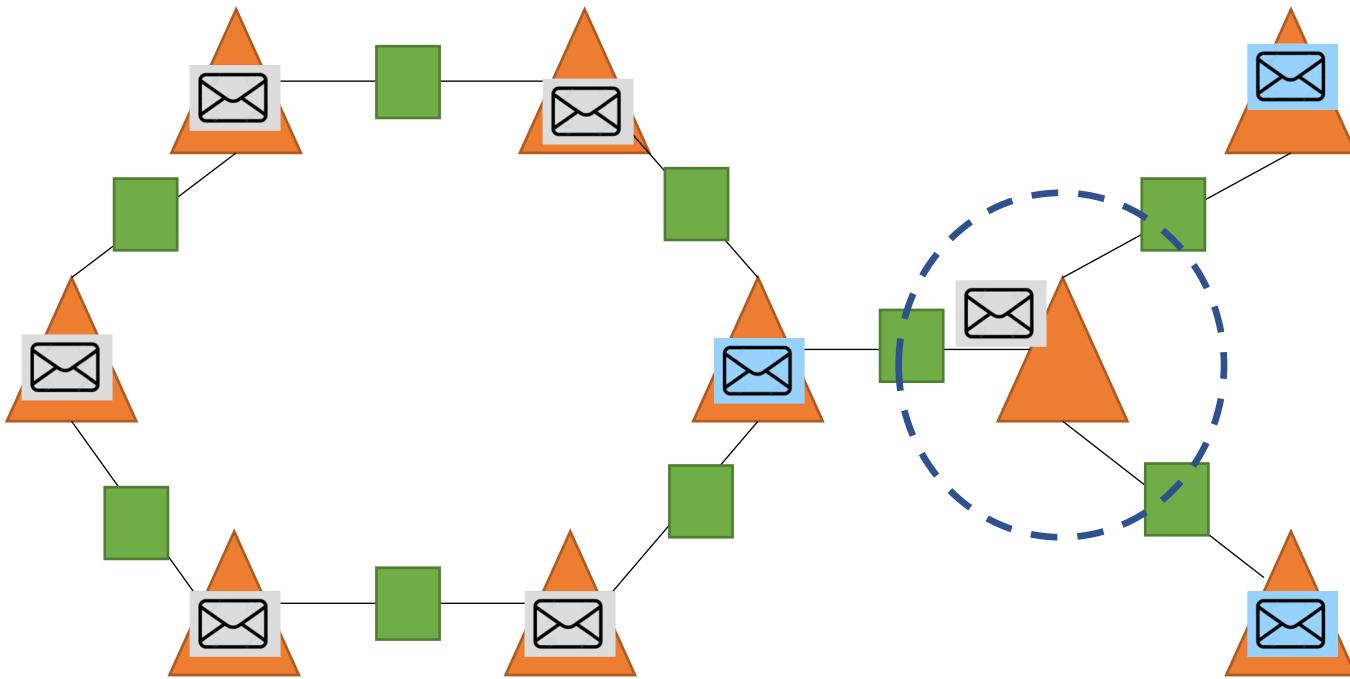












$$\text{✉} = \Delta(\text{✉}, \sum \text{✉})$$

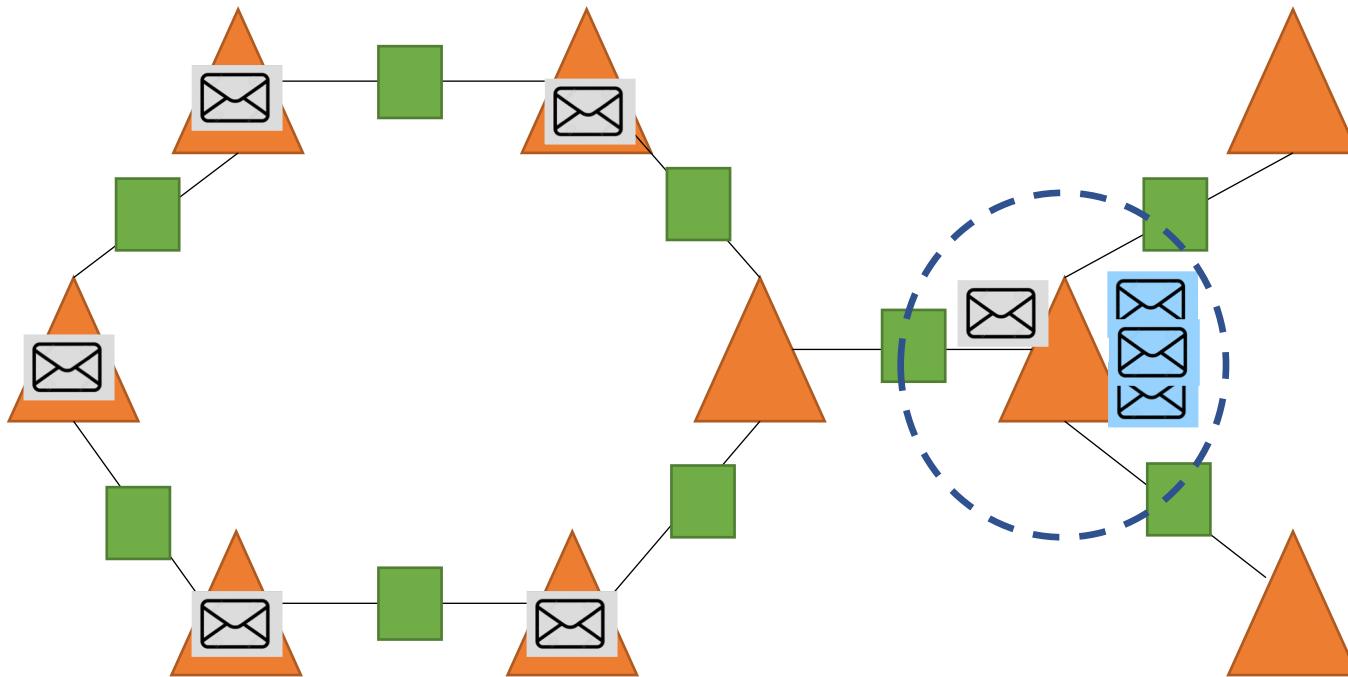


TABLE 2
Different variants of graph neural networks.

Name	Variant	Aggregator	Updater
Spectral Methods	ChebNet	$\mathbf{N}_k = \mathbf{T}_k(\tilde{\mathbf{L}})\mathbf{X}$	$\mathbf{H} = \sum_{k=0}^K \mathbf{N}_k \Theta_k$
	1 st -order model	$\mathbf{N}_0 = \mathbf{X}$ $\mathbf{N}_1 = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = \mathbf{N}_0 \Theta_0 + \mathbf{N}_1 \Theta_1$
	Single parameter	$\mathbf{N} = (\mathbf{I}_N + \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}) \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$
	GCN	$\mathbf{N} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$
Non-spectral Methods	Neural FPs	$\mathbf{h}'_{\mathcal{N}_v} = \mathbf{h}_v^{t-1} + \sum_{k=1}^{\mathcal{N}_v} \mathbf{h}_k^{t-1}$	$\mathbf{h}_v^t = \sigma(\mathbf{h}'_{\mathcal{N}_v} \mathbf{W}_L^{\mathcal{N}_v})$
	DCNN	Node classification: $\mathbf{N} = \mathbf{P}^* \mathbf{X}$ Graph classification: $\mathbf{N} = \mathbf{1}_N^T \mathbf{P}^* \mathbf{X} / N$	$\mathbf{H} = f(\mathbf{W}^e \odot \mathbf{N})$
	GraphSAGE	$\mathbf{h}'_{\mathcal{N}_v} = \text{AGGREGATE}_v(\{\mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v\})$	$\mathbf{h}_v^t = \sigma(\mathbf{W}^t \cdot [\mathbf{h}_v^{t-1} \ \mathbf{h}'_{\mathcal{N}_v}])$
Graph Attention Networks	GAT	$\alpha_{ek} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_e \mathbf{W}\mathbf{h}_k]))}{\sum_{j \in \mathcal{N}_v} \exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_e \mathbf{W}\mathbf{h}_j]))}$ $\mathbf{h}'_{\mathcal{N}_v} = \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{ek} \mathbf{W}\mathbf{h}_k)$ Multi-head concatenation: $\mathbf{h}'_{\mathcal{N}_v} = \left\ _{m=1}^M \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{ek}^m \mathbf{W}^m \mathbf{h}_k)$ Multi-head average: $\mathbf{h}'_{\mathcal{N}_v} = \sigma\left(\frac{1}{M} \sum_{m=1}^M \sum_{k \in \mathcal{N}_v} \alpha_{ek}^m \mathbf{W}^m \mathbf{h}_k\right)$	$\mathbf{h}_v^t = \mathbf{h}'_{\mathcal{N}_v}$
Gated Graph Neural Networks	GGNN	$\mathbf{h}'_{\mathcal{N}_v} = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1} + \mathbf{b}$	$\mathbf{z}_v^t = \sigma(\mathbf{W}^z \mathbf{h}_v^t + \mathbf{U}^z \mathbf{h}_v^{t-1})$ $\mathbf{r}_v^t = \sigma(\mathbf{W}^r \mathbf{h}_v^t + \mathbf{U}^r \mathbf{h}_v^{t-1})$ $\tilde{\mathbf{h}}_v^t = \tanh(\mathbf{W}^h \mathbf{h}_v^t + \mathbf{U}(\mathbf{r}_v^t \odot \mathbf{h}_v^{t-1}))$ $\mathbf{h}_v^t = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{t-1} + \mathbf{z}_v^t \odot \tilde{\mathbf{h}}_v^t$
Graph LSTM	Tree LSTM (Child sum)	$\mathbf{h}'_{\mathcal{N}_v} = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{U}^i \mathbf{h}_{\mathcal{N}_v}^{t-1} + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{U}^f \mathbf{h}_{vk}^{t-1} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{U}^o \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{U}^u \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{k \in \mathcal{N}_v} \mathbf{f}_{vk}^t \odot \mathbf{c}_k^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$
	Tree LSTM (N-ary)	$\mathbf{h}_{\mathcal{N}_v}^{ti} = \sum_{l=1}^K \mathbf{U}_l^i \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v k}^{tf} = \sum_{l=1}^K \mathbf{U}_{kl}^f \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{to} = \sum_{l=1}^K \mathbf{U}_l^o \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{tu} = \sum_{l=1}^K \mathbf{U}_l^u \mathbf{h}_{vl}^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{ti} + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v k}^{tf} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{to} + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{tu} + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{l=1}^K \mathbf{f}_{vl}^t \odot \mathbf{c}_l^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$
	Graph LSTM in [44]	$\mathbf{h}_{\mathcal{N}_v}^{ti} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^i \mathbf{h}_k^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{to} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^o \mathbf{h}_k^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{tu} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^u \mathbf{h}_k^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{ti} + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{U}_{m(v,k)}^f \mathbf{h}_k^{t-1} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{to} + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{tu} + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{k \in \mathcal{N}_v} \mathbf{f}_{vk}^t \odot \mathbf{c}_k^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$

Two view of graph neural nets

- Spectral networks
 - Rely on the eigen-decomposition of the Laplacian matrix
 - Any perturbation to a graph results in a change of eigenbasis
 - Need to load the whole graph into the memory to perform graph convolution
- Spatial networks
 - Take the aggregation of the central node representation and its neighbors representation

Spectral networks

- What the heck is “spectral” now?
 - It simply means *decomposing* a signal/audio/image/graph into a combination (usually, a sum) **of simple elements** (wavelets, graphlets)
 - these simple elements are usually orthogonal, i.e. mutually linearly independent, and therefore form a basis

Spectral networks

- “Spectral” in signal processing
 - Usually means applying the Fourier Transform to a signal
 - Decompose it into elementary sine and cosine waves of different frequencies
- “Spectral” in Graph signal processing
 - Eigen-decomposition of the graph Laplacian L
 - Decompose it into communities

Spectral Networks

- Given a graph, $G = (V, E)$
 - X , a feature description for every node ($N * D$)
 - Adjacency matrix (A)
- Output Z , node level output ($N * F$)

Spectral Networks

- Every neural network layer can be written as

$$H^{(l+1)} = f(H^{(l)}, A)$$

$$H^{(0)} = X$$

- For example:

$$f(H^{(l)}, A) = \sigma(AH^{(l)}W^{(l)})$$

- For every node, sum up all the feature vectors of all neighboring nodes but not the node itself

Spectral Networks

- Fix it by adding the identity matrix

$$\hat{A} = A + I$$

- If we want to add multiple layers to one another, they have to have similar scales
- Need to normalize the adjacency matrix, such that all rows of the adjacency matrix must add up to one

$$\hat{D}^{-0.5} \hat{A} \hat{D}^{-0.5}$$

- \hat{D} is the diagonal node degree matrix of \hat{A}

Spectral Networks

- Putting it all together

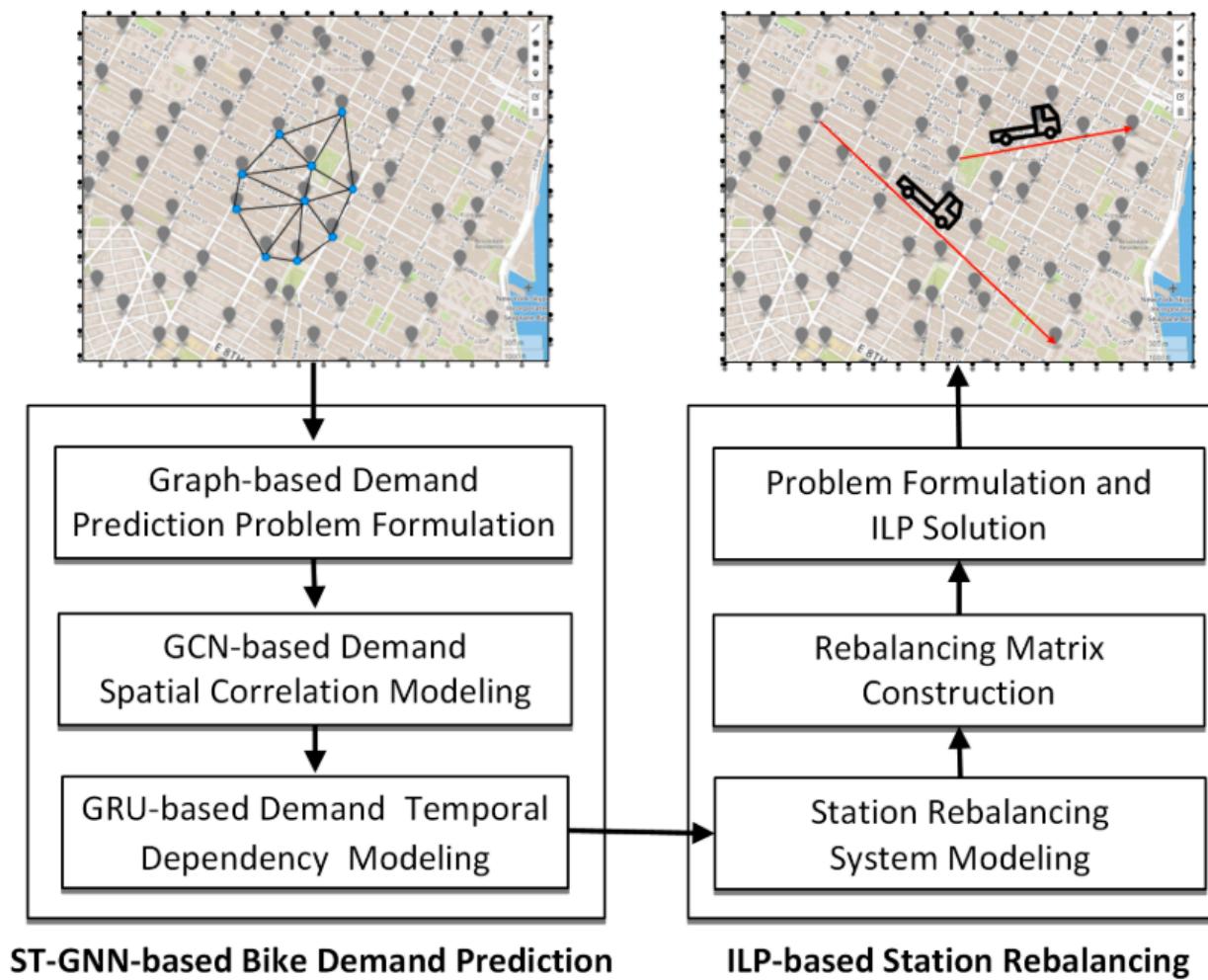
$$f(H^{(l)}, A) = \sigma(\hat{D}^{-0.5} \hat{A} \hat{D}^{-0.5} H^{(l)} W^{(l)})$$

- Compare with:

$$f(H^{(l)}, A) = \sigma(A H^{(l)} W^{(l)})$$

BikeNet: Accurate Bike Demand Prediction Using Graph Neural Networks for Station Rebalancing

Ruiying Guo¹, Zhihan Jiang¹, Jingchun Huang¹, Jianrong Tao², Cheng Wang¹, Jonathan Li^{1,3}, Longbiao Chen*¹



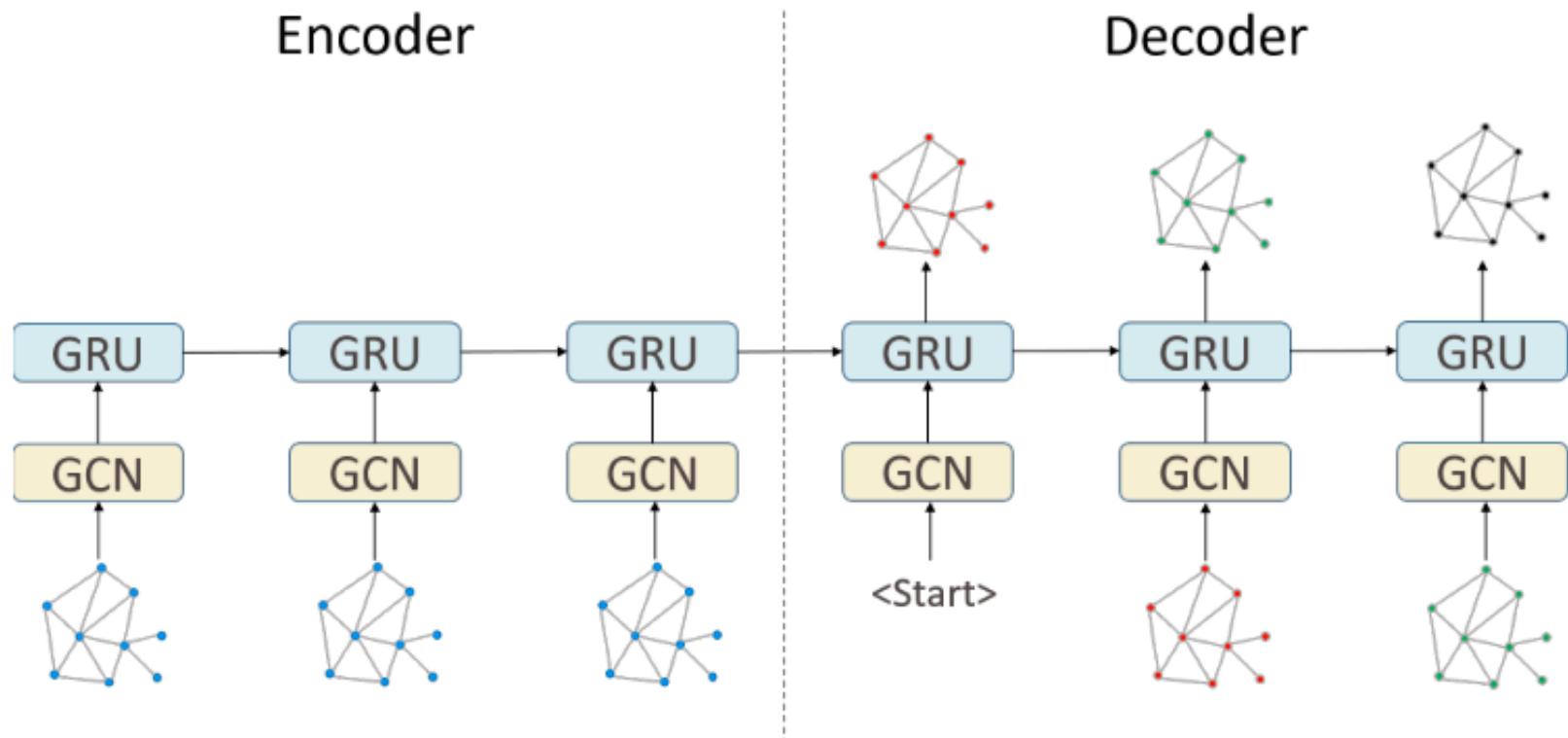
Demand prediction for rebalancing



Demand prediction for rebalancing

- Represent the bike sharing system as a weighted undirected spatiotemporal graph
- Node v 's message is bike demand v during $[t, t + \Delta t]$
- Edges: distance between stations
 - Could use correlations
 - Mutual information
 - Etc.

Demand prediction for rebalancing



Demand prediction for rebalancing

$$\text{minimize} \quad C^T x \quad (3)$$

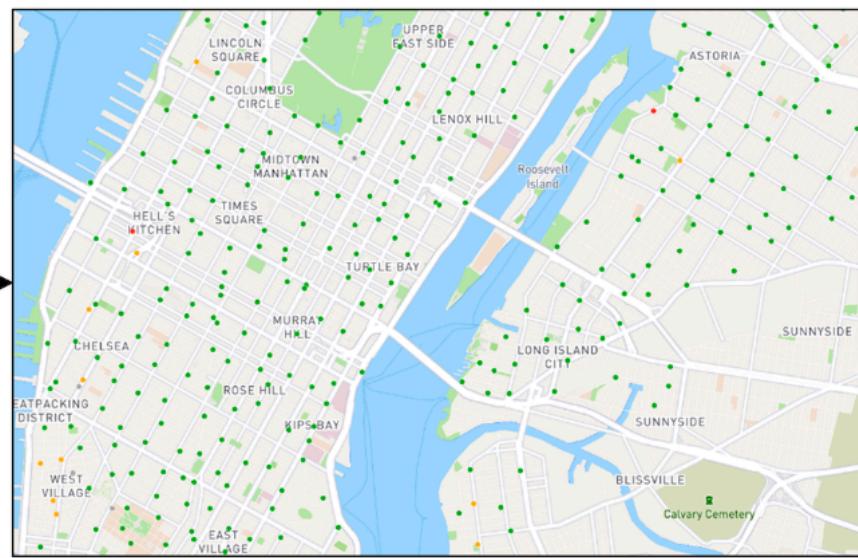
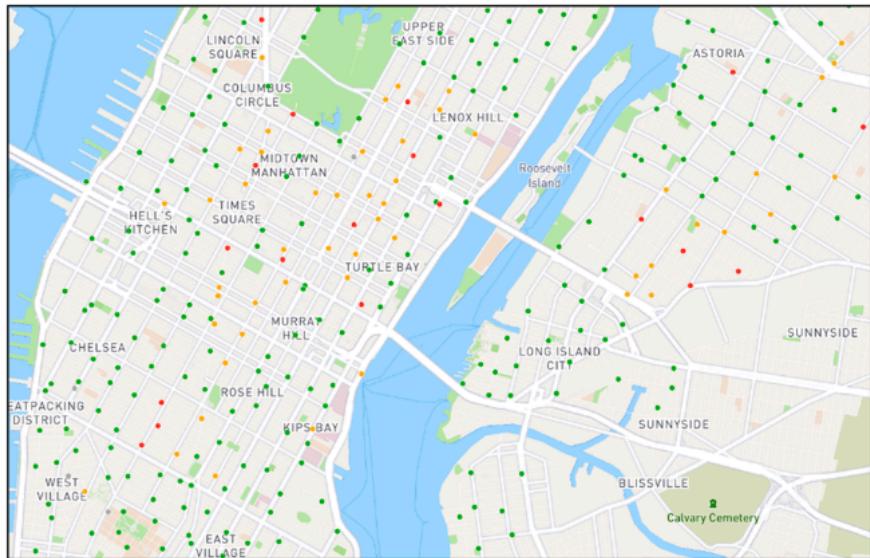
subject to

$$\sum_{j=1}^l a_{ij}x_j \leq 0 \quad (i = 1, \dots, n+m) \quad (4)$$

$$x_j \quad (j = 1, \dots, n+m) \geq 0 \quad (5)$$

$$\sum_{i=1}^{n+m} \sum_{j=1}^l a_{ij}x_j = 2\min(n, m) \quad (6)$$

Demand prediction for rebalancing



References

CMU Deep learning course 2019
Bhiksha Raj