

Personal Recommendation Using Deep Recurrent Neural Networks in NetEase

2016.08.29

서지혜

Personal Recommendation Using Deep Recurrent Neural Networks in NetEase

Sai Wu ^{#1}, Weichao Ren ^{#2}, Chengchao Yu ^{#3}, Gang Chen ^{#4}, Dongxiang Zhang ^{♢5}, Jingbo Zhu ^{◇6}

[#] *College of Computer Science and Technology, Zhejiang University, Hangzhou, China*

^{1,2,3,4} {wusai, weichaor, 21311200, cg}@zju.edu.cn

[♢] *School of Computer Science and Engineering, University of Electronic Science and Technology of China*

⁵ zhangdongxiang37@gmail.com

[◇] *NetEase (Hangzhou) Network Co., Ltd., Hangzhou, China*

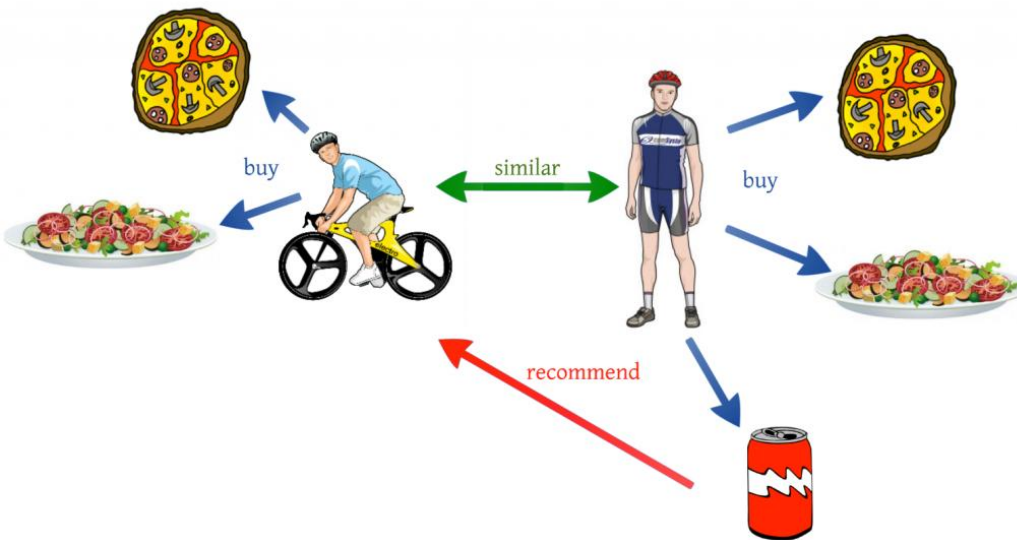
⁶ zhujingbo@corp.netease.com

Abstract—Each user session in an e-commerce system can be modeled as a sequence of web pages, indicating how the user interacts with the system and makes his/her purchase. A typical recommendation approach, e.g., Collaborative Filtering, generates its results at the beginning of each session, listing the most likely purchased items. However, such approach fails to exploit current viewing history of the user and hence, is unable to provide a real-time customized recommendation service. In this paper, we build a deep recurrent neural network to address

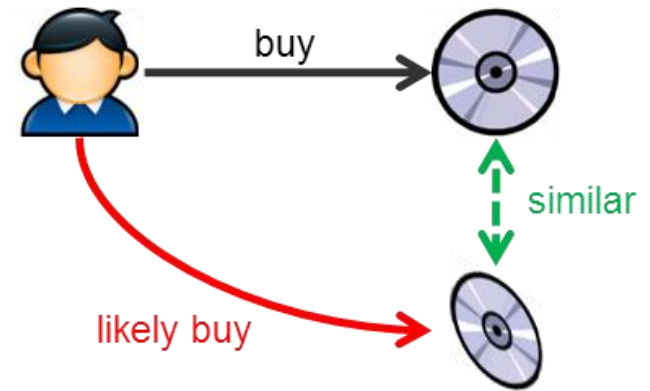
session, which update their values throughout the session. Before making his/her purchase, a user will view a long list of web pages (In Kaola, a user views 12 pages averagely before adding one item to the cart). This viewing history actually records how users interact with the system and perform their purchases, while previous CF recommendation approaches fail to exploit it. In Figure 1, the first page of the session is “index.html”. As no specific information can be retrieved, we will return our recommendations using the CF algorithm as the

Collaborative Filtering


User-based Recommendation



Item-based Recommendation



Personal Recommendation Using Deep Recurrent Neural Networks in NetEase

- ▶ Collaborative Filtering
 - Unable to provide a real-time customized recommendation service
 - Established using stale data
 - ▶ In this paper
 - Using a history state
 - Build a DRNN with FNN
 - Develop optimizer to automatically tune the parameters of NN
- 

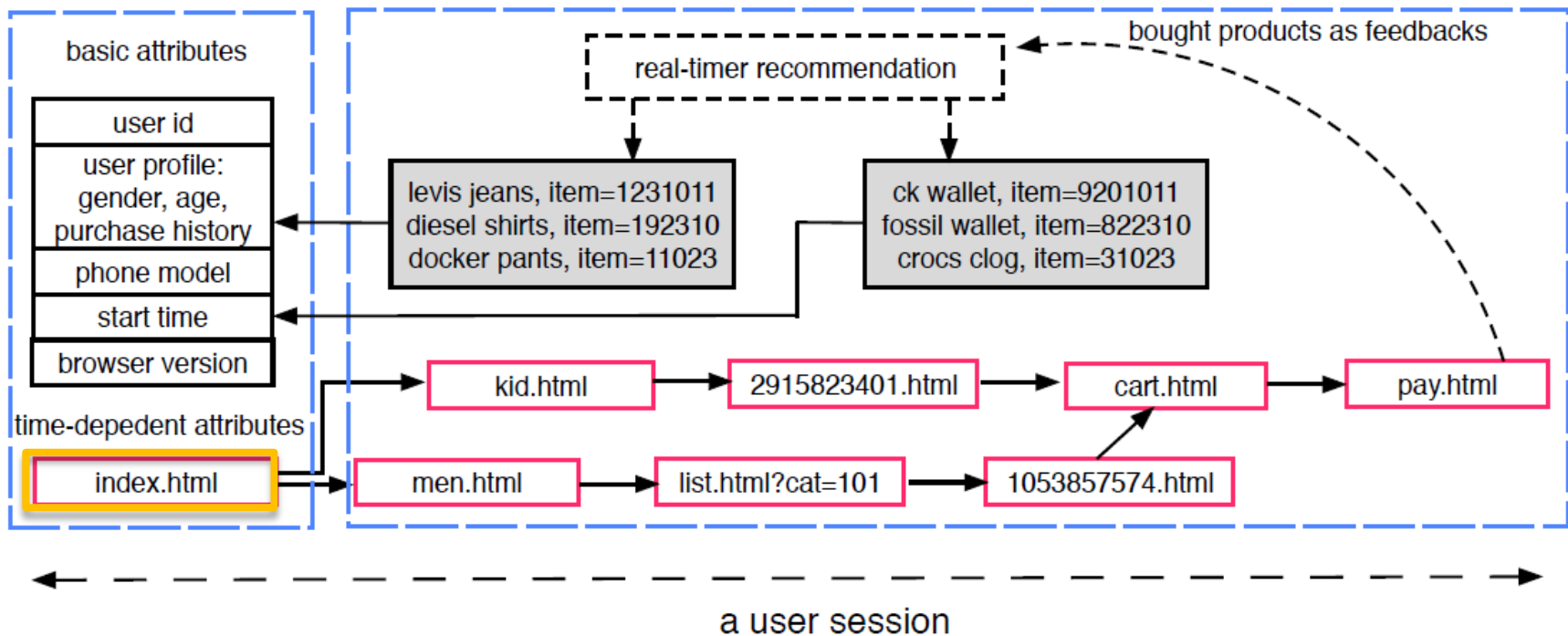
Real-time Recommendation in E-commerce System



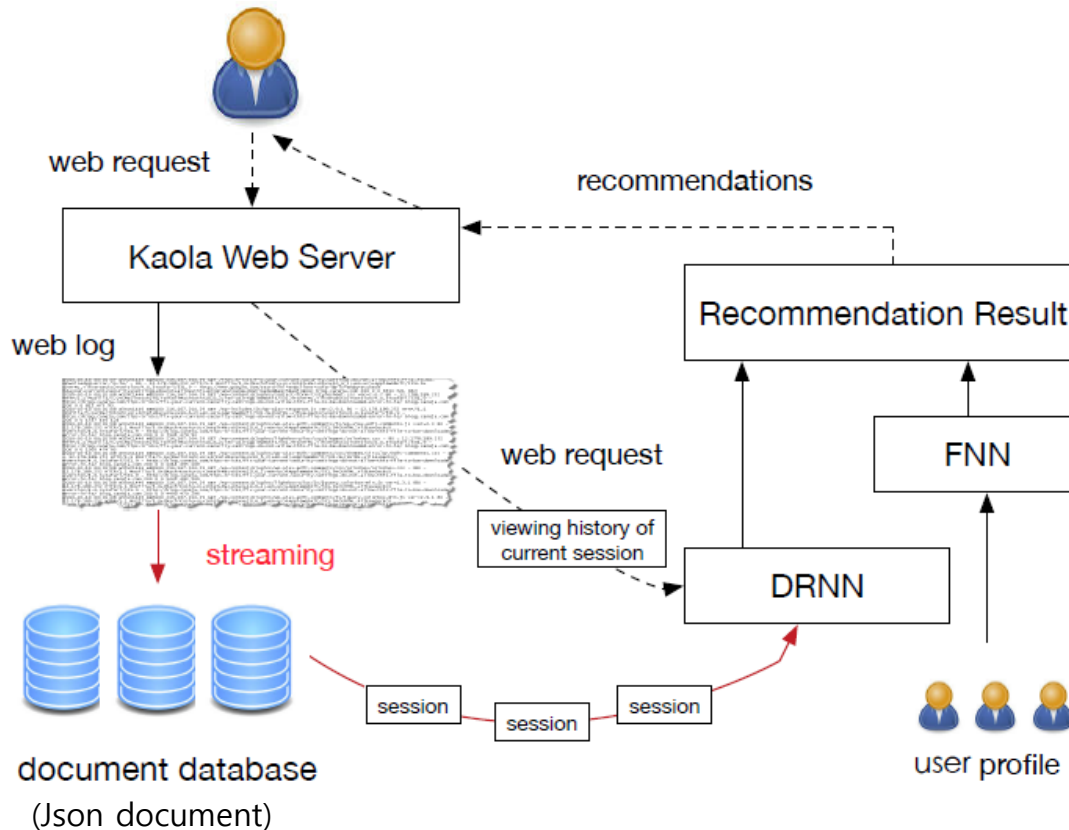
CUSTOMERS

Men's
wallets &
kids' shoes

e-commerce system



Overview of Recommendation Module

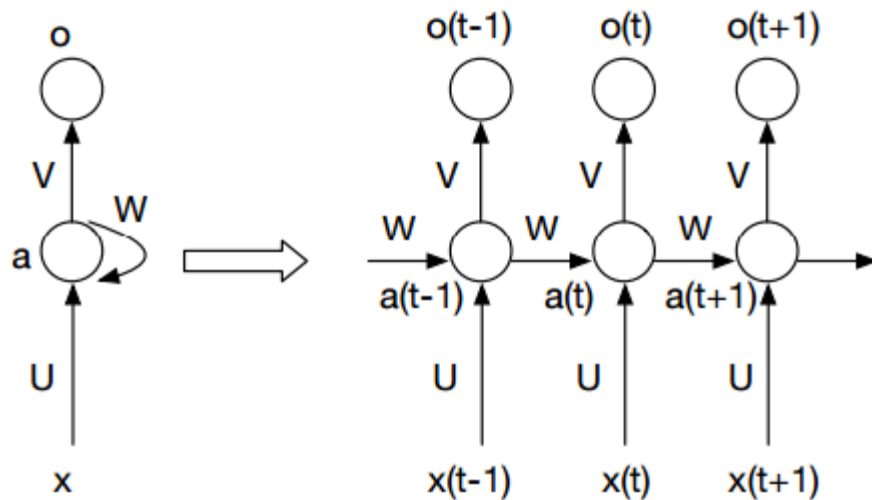


- ▶ Consider both Interest of current session (DRNN) and past interest (FNN)
- ▶ Session continues Gradually refine
- ▶ Session completes Adjust our model using new training sample

Fig. 2. Work Flow of Recommendation Module

Deep Recurrent Neural Network

► RNN



$$a(i) = f(Ux(i) + Wa(i-1))$$

$$o(i) = \text{softmax}(Va(i))$$

Deep Recurrent Neural Network

► Deep RNN Model with Infinite State

Input $\rightarrow \{p_0, p_1, \dots, p_{n-1}\}$ (a sequence of web pages)

Output $\rightarrow \{v_0, v_1, \dots, v_{N-1}\}$

◦ Update function

$a_i(t)$ denote the state of the i th layer at state t .

$$a_i(t) = \begin{cases} f(W_i a_i(t-1) + Z_i(a_{i-1}(t) + b_i(t))) & \text{if } i > 1 \\ f(W_i a_i(t-1) + Z_i(V_t + \theta(p_t))) & \text{if } i = 1 \end{cases}$$

W_i : weights of connections from state $t-1$

Z_i : weights of connections from the lower hidden layer at the same state

$b_i(t)$: bias

$\theta(p_t)$: bias for input layer

Deep Recurrent Neural Network

▶ Deep RNN Model with Infinite State

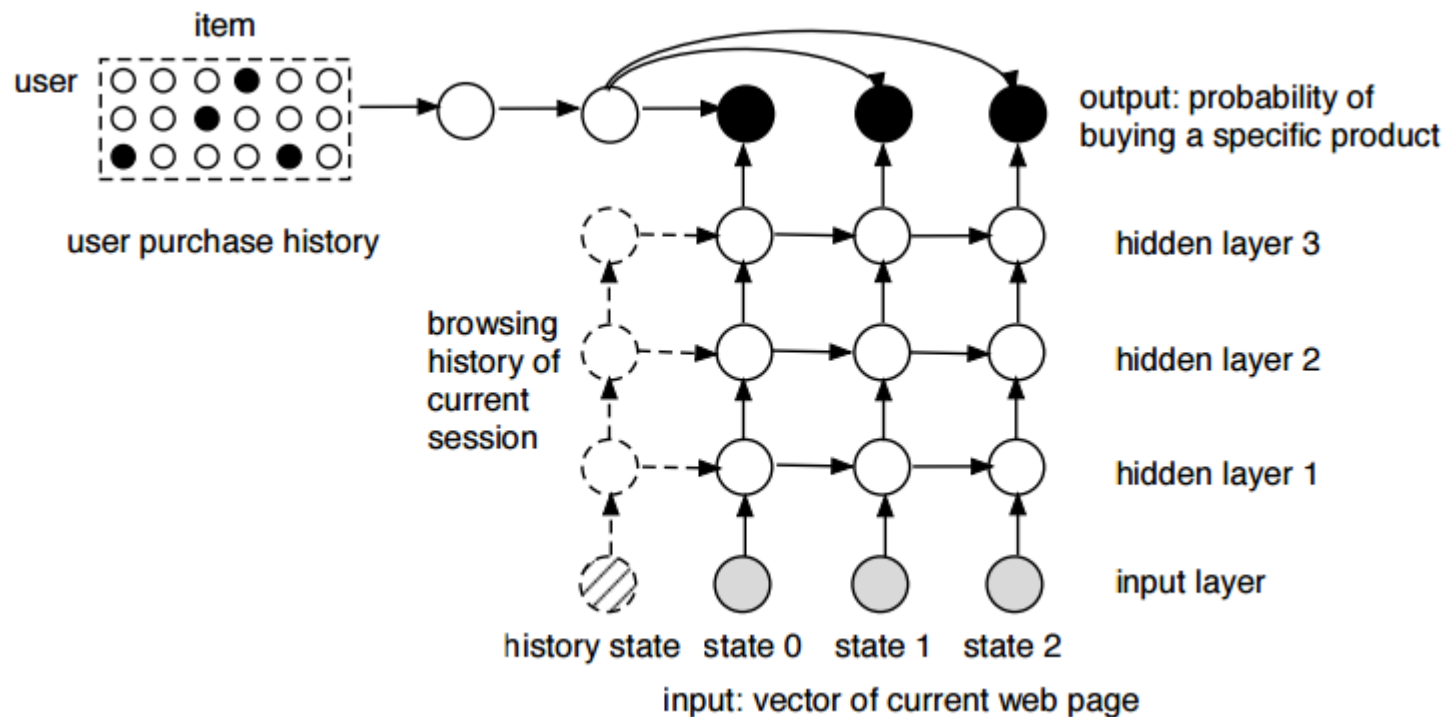


Fig. 4. Illustration of Deep RNN model

Deep Recurrent Neural Network

▶ Deep RNN Model with History State

- All pages : $\{p_0, p_1, \dots, p_x\}$ If $x \leq n \rightarrow$ history state is not used.
- Input for the history state

$$\bar{V} = \sum_{i=0}^{x-n} \epsilon_i V_i$$

V_i : vector for page p_i
 ϵ_i : the aging factor for old state

$$\epsilon_i = \frac{\theta(p_i)}{\sum_{j=i}^{x-n} \theta(p_j)}$$

Deep Recurrent Neural Network

- ▶ Collaboration with Collaborative Filtering
 - CF (Feedforward Neural Network)
 - Generates a good recommendation if users follow their old purchase patterns
 - RNN
 - Effectively predict the unexpected purchase of a specific user

Deep Recurrent Neural Network

► Collaboration with Collaborative Filtering

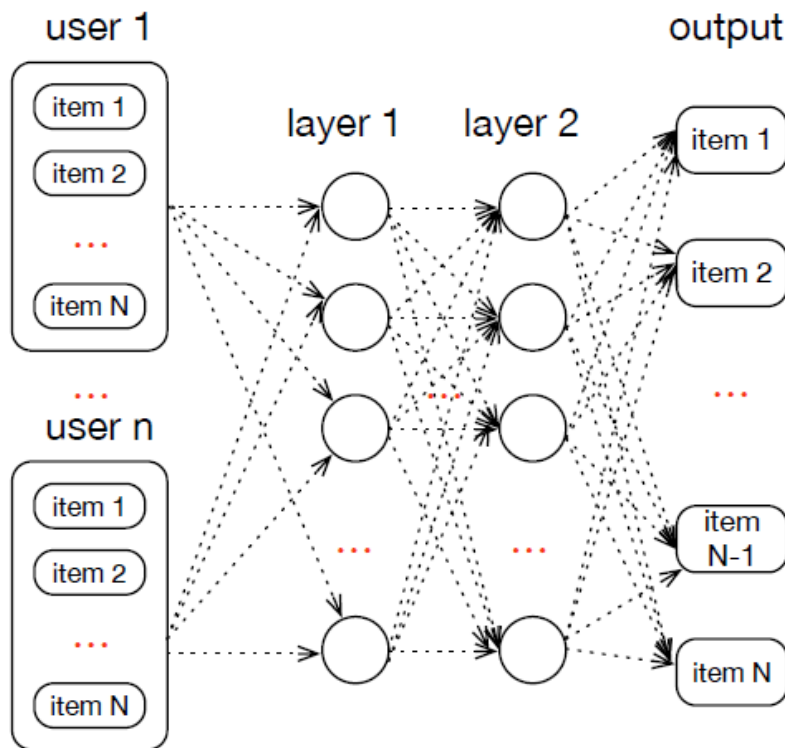


Fig. 5. FNN for CF

- State update

$$\bar{a}_j^{(i)} = f\left(\sum_{x=0}^{\bar{E}-1} w_j^{(i-1)} \bar{a}_x^{(i-1)} + b_x^{(i-1)}\right)$$

$\bar{a}_j^{(i)}$: jth neurons at the ith hidden layer

\bar{E} : number of each hidden layer neurons

$w_j^{(i-1)}$: weight from the last layer

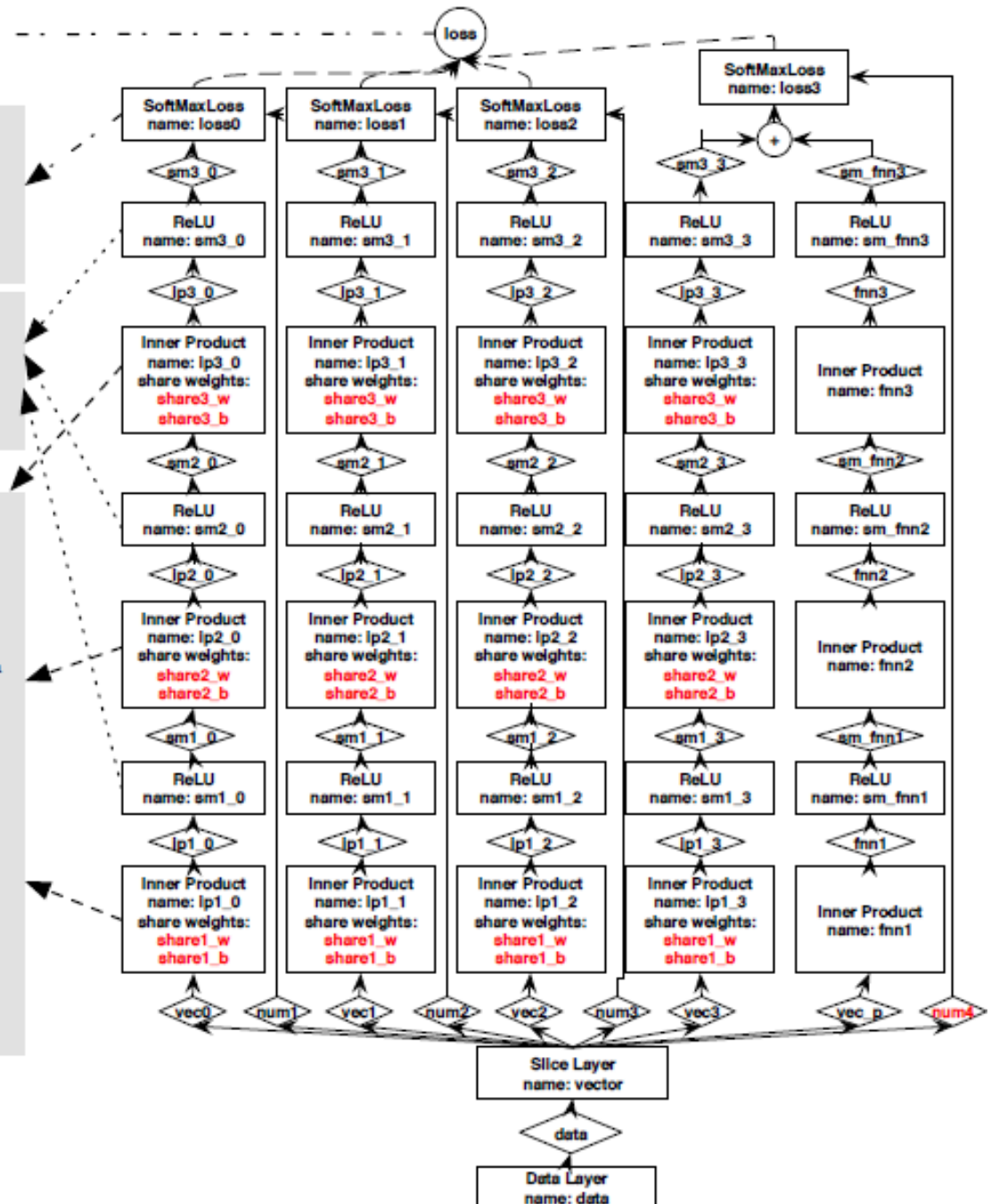
$b_x^{(i-1)}$: corresponding bias

DRNN Implementations

```
layer {
  name: "loss0"           # layer name
  type: "SoftmaxWithLoss" # layer type
  bottom: "sm3_0"         # 1st input (a vector)
  bottom: "label0"        # 2nd input (a scalar)
  top: "loss0"            # loss output blob
}
```

```
layer {
  name: "sm1_0"           # layer name
  type: "ReLU"            # layer type
  bottom: "ip1_0"          # input blob
  top: "sm1_0"            # output blob
}
```

```
layer {
  name: "ip1_0"           # layer name
  type: "InnerProduct"    # layer type
  bottom: "vec0"          # input blob
  top: "ip1_0"            # output blob
  inner_product_param {
    num_output: 1000
    weight_filler {
      # Matrix W initialization
      type: "gaussian"
      std: 1
    }
    bias_filler {
      # bias vector values
      type: "constant"
      value: 1
    }
  }
  param {
    name: "share1_w"       # share matrix W values
  }
  param {
    name: "share1_b"       # share bias values name
  }
}
```



Automatic Optimization Framework

▶ Automatic Code Generation

- Parameters : number of states/neurons, type of activation functions and the loss function
- Code generator
 - Accepts a configuration file defining the values of parameters
 - Outputs the Caffe script for the corresponding DRNN model
- Basic parameters / network structure parameters

Automatic Optimization Framework

► Automatic Code Generation

Algorithm 1 CodeGen(int w , int l , int h)

```
1: Layer[] input = genInputLayer()
2: for  $i = 0$  to  $w - 1$  do
3:   Layer top = null, Layer btm = null
4:   for  $j = 0$  to  $l - 1$  do
5:     if btm  $\neq$  null then
6:       btm = genNeuronLayer()
7:       top = genNeuronLayer()
8:       connectLayer(btm, top)
9:       btm = top
10:  Layer loss = genLossLayer()
11:  connectLayer(btm, loss)
12: Layer history_top = null, Layer history_btm = null
13: for  $j = 0$  to  $l - 1$  do
14:   if history_btm  $\neq$  null then
15:     btm = genNeuronLayer()
16:     history_top = genNeuronLayer()
17:     connectLayer(history_btm, history_top)
18:     history_btm = history_top
19: for  $i = 0$  to  $w - 1$  do
20:   Layer loss = getLossLayer( $i$ )
21:   Layer next = input[ $j$ ].getTopLayer()
22:   for  $j = 0$  to  $h - 1$  do
23:     connectLayer(next, loss)
24:     next = next.getTopLayer()
```

Creates all neuron layers
for current states

Sets up the history state
without the loss layer

Connect neuron layers to
their loss layers to create
the shortcuts

Automatic Optimization Framework

► Model Tuning

Algorithm 2 GenTune(int w , int l , int h)

```
1:  $\tau = 0$ 
2:  $P(\tau) = \text{GeneratePopulation}()$ 
3:  $\text{maxFit} = \text{maximum fitness of chromosomes in } P(\tau)$ 
4:  $\text{best} = \text{chromosome with maximum fitness in } P(\tau)$ 
5: while  $\tau < N$  and user does not interrupt do
6:    $\text{randomShuffle}(P(\tau))$ 
7:   for  $i = 1$  to  $\frac{P(\tau).size}{2}$  do
8:      $c_1 = P(\tau)[i * 2 - 1]$ 
9:      $c_2 = P(\tau)[i * 2]$ 
10:     $r_1, r_2$  and  $r_3$  are random values in  $[0, 1]$ 
11:    if  $r_1 < p_c$  then
12:       $\text{crossover}(c_1, c_2)$ 
13:    if  $r_2 < p_m$  then
14:       $\text{mutate}(c_1)$ 
15:    if  $r_3 < p_m$  then
16:       $\text{mutate}(c_2)$ 
17:    if  $r_1 < p_c$  or  $r_2 < p_m$  then
18:       $\text{maxFit} = \text{Max}(\text{getFitness}(c_1), \text{maxFit})$ 
19:       $\text{best} = \text{chromosome with maximum fitness in new } P(\tau)$ 
20:    if  $r_1 < p_c$  or  $r_3 < p_m$  then
21:       $\text{maxFit} = \text{Max}(\text{getFitness}(c_2), \text{maxFit})$ 
22:       $\text{best} = \text{chromosome with maximum fitness in new } P(\tau)$ 
23:     $\tau++$ 
```

$P(\tau)$: chromosome set in the τ th iteration

$$fit = accuracy + \frac{1}{1 + loss}$$

Experiments

► Data

- Web log : 232,326 records (June 1st , 2015)
- 27,985 session 37,667 unique users
- 60% training , 20% validation, 20% testing

► Equipment

- Xeon 2.6G 8-core CPUs / 64GB memory
- One GeForce GTX Titian Z GPU (12GB DDR5)
- Caffe (GPU mode)

► DRNN model structure

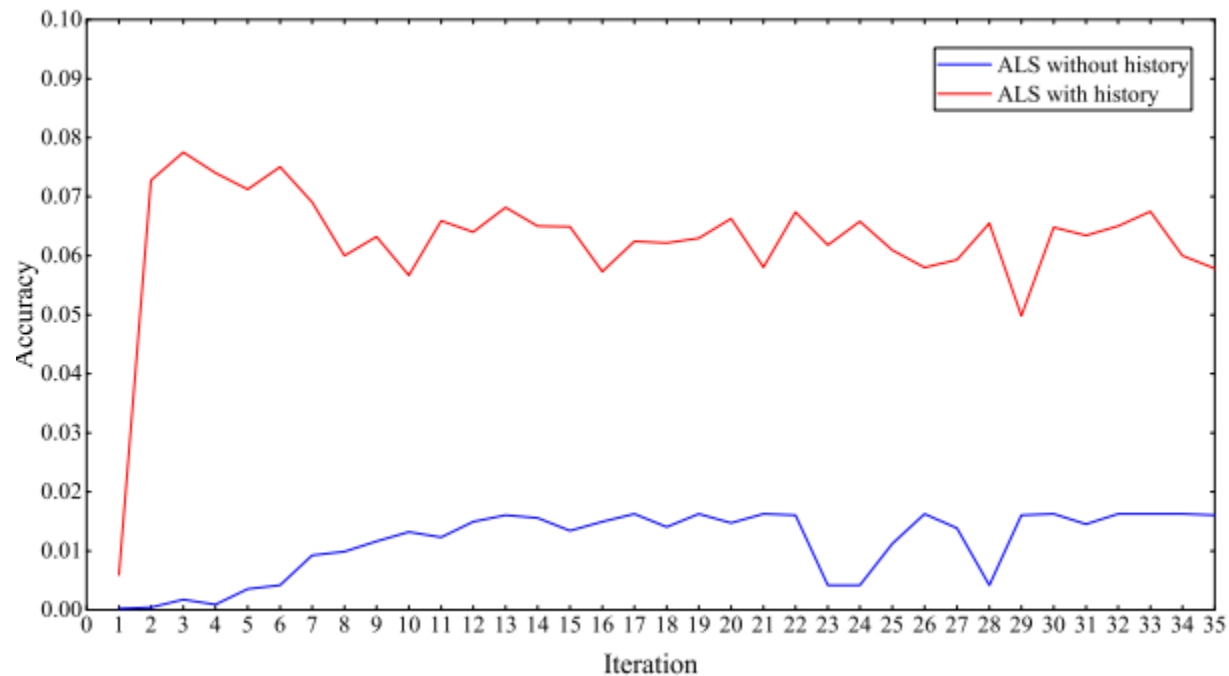
- 4 state
 - One history state
- 3 hidden layers

TABLE I. EXPERIMENT SETTINGS

Parameter	Default Value
Number of Hidden Layers in RNN	3
Number of Hidden Layers in FNN	3
Number of States in RNN	4
Number of Neurons in RNN	1000,100
Number of Neurons in FNN	1000
Activation Function	ReLU
Loss Function	Softmax
Initial Learning Rate	0.1
Batch Size in RNN	9000
Batch Size in FNN	5000
Weight Initialization Function	Xavier

Experiments

► Performance of ALS



Experiments : Effect of Batching

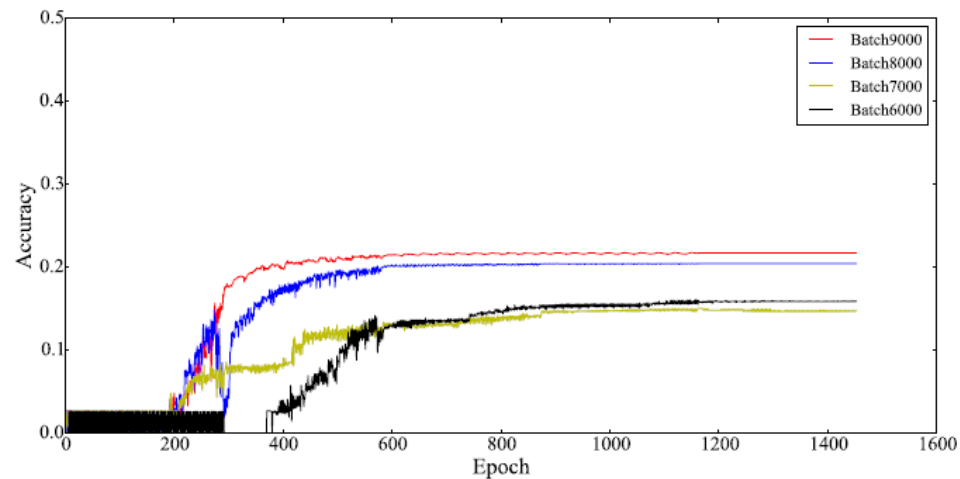


Fig. 9. Effect of Batching without Tuning

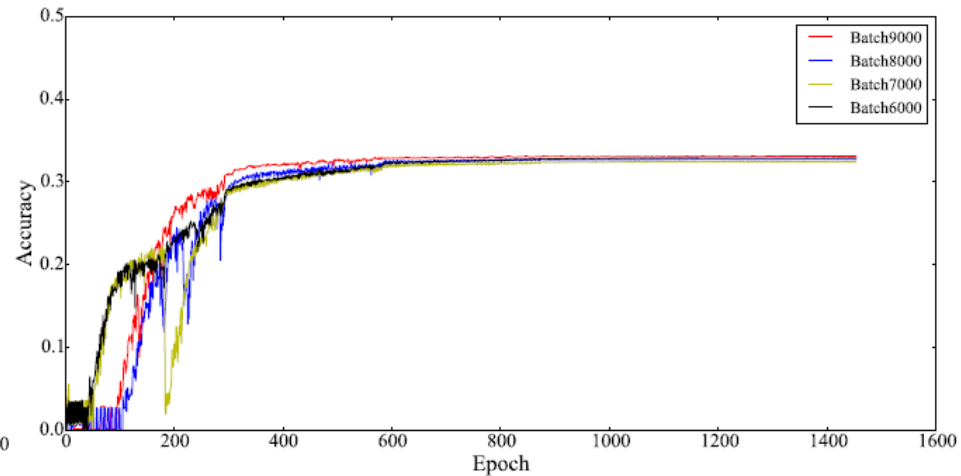


Fig. 10. Effect of Batching with Tuning

Experiments : Effect of FNN

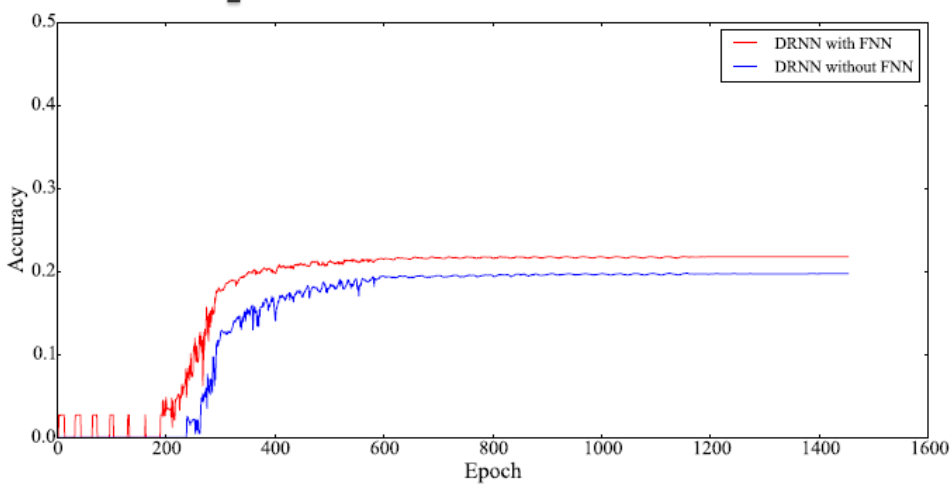


Fig. 11. Effect of FNN without Tuning

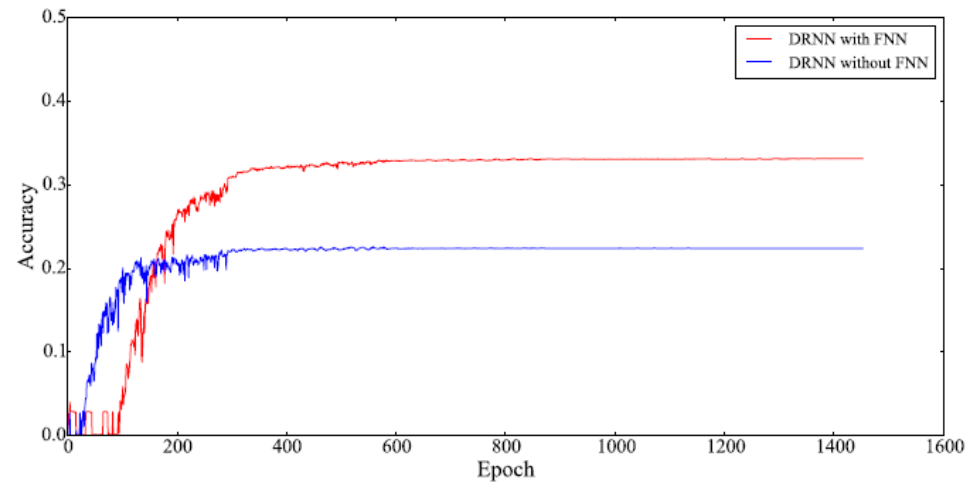


Fig. 12. Effect of FNN with Tuning

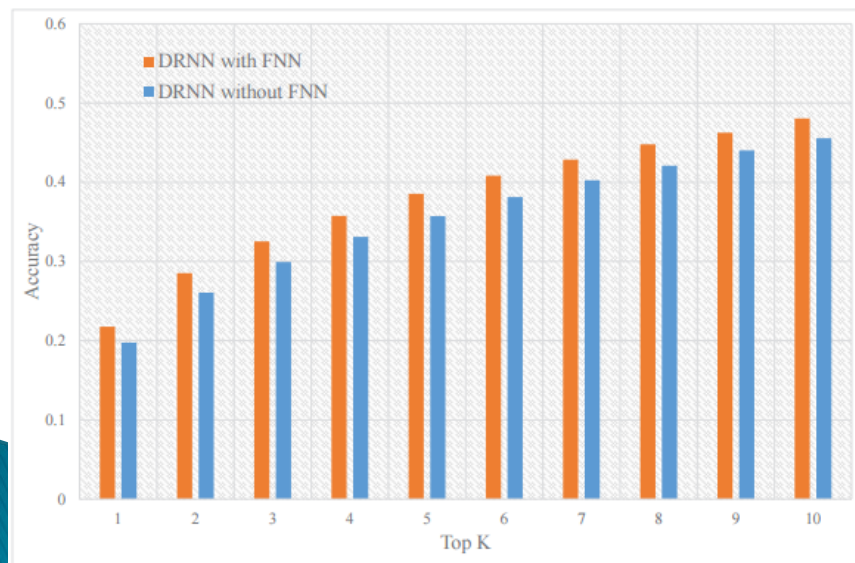


Fig. 13. FNN Top-K Result without Tuning

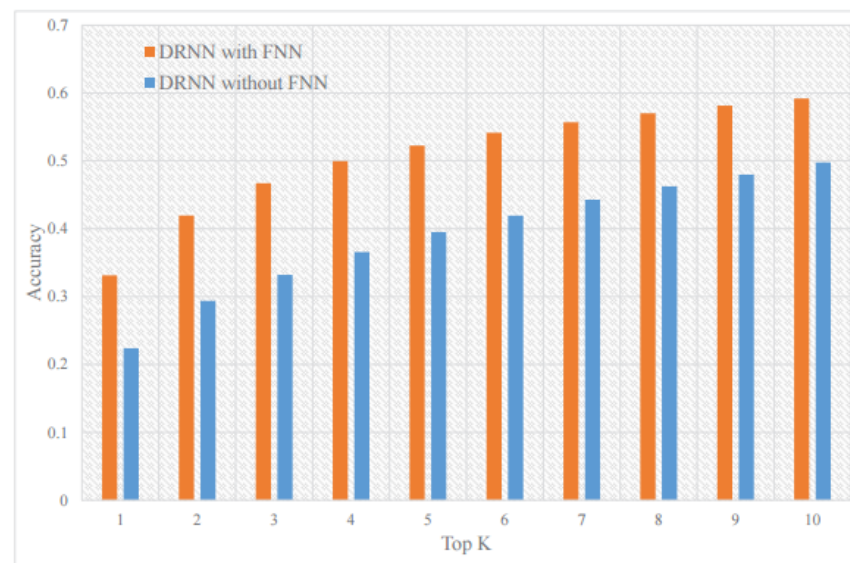


Fig. 14. FNN Top-K Result with Tuning

Experiments : Effect of History state

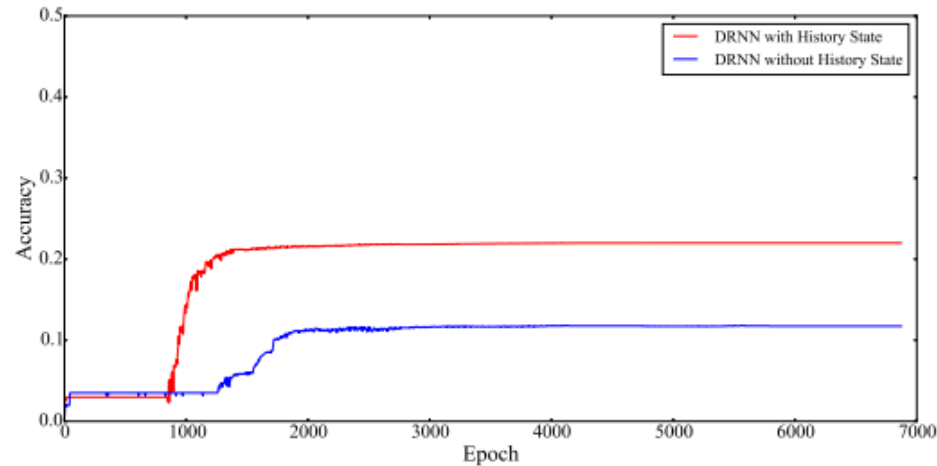


Fig. 15. Effect of History State without Tuning

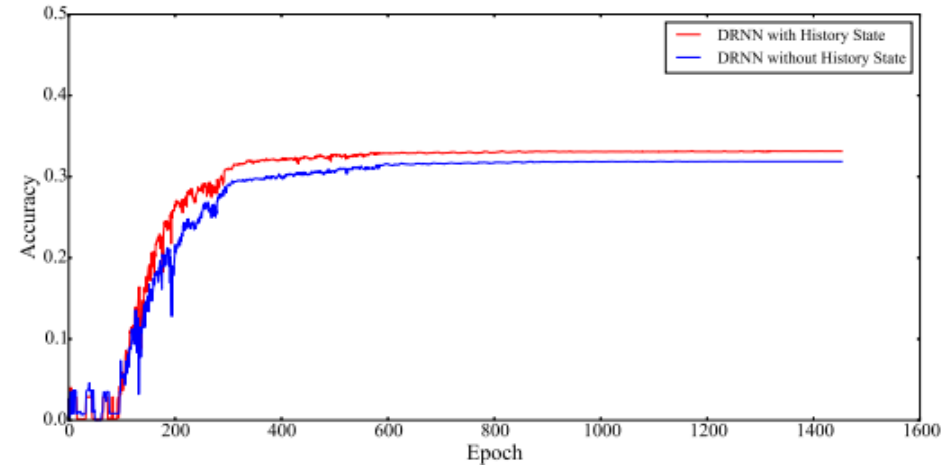


Fig. 16. Effect of History State with Tuning



Fig. 17. History State Top-K Result without Tuning

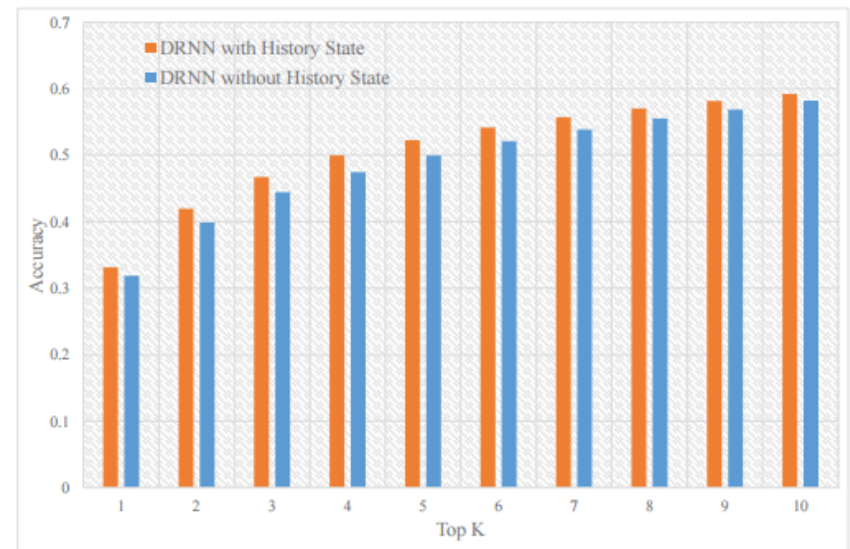


Fig. 18. History State Top-K Result with Tuning

Experiments

► Convergence Rate

- Initial learning rate

$$\phi = 0.1$$

- Every 100 epochs

$$\phi = \frac{\phi}{5}$$

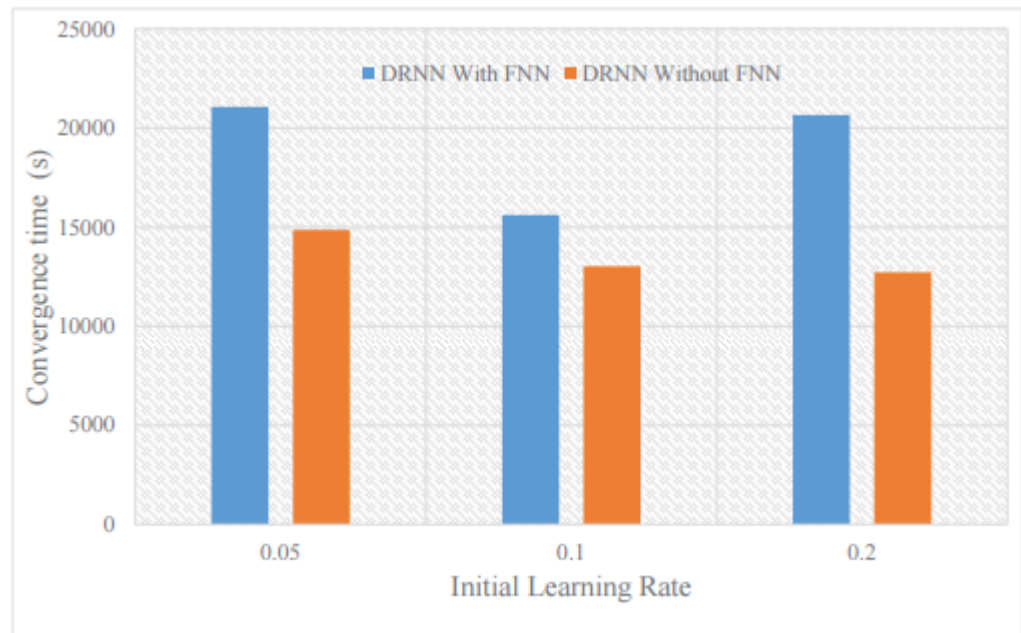


Fig. 19. Convergence Rate