

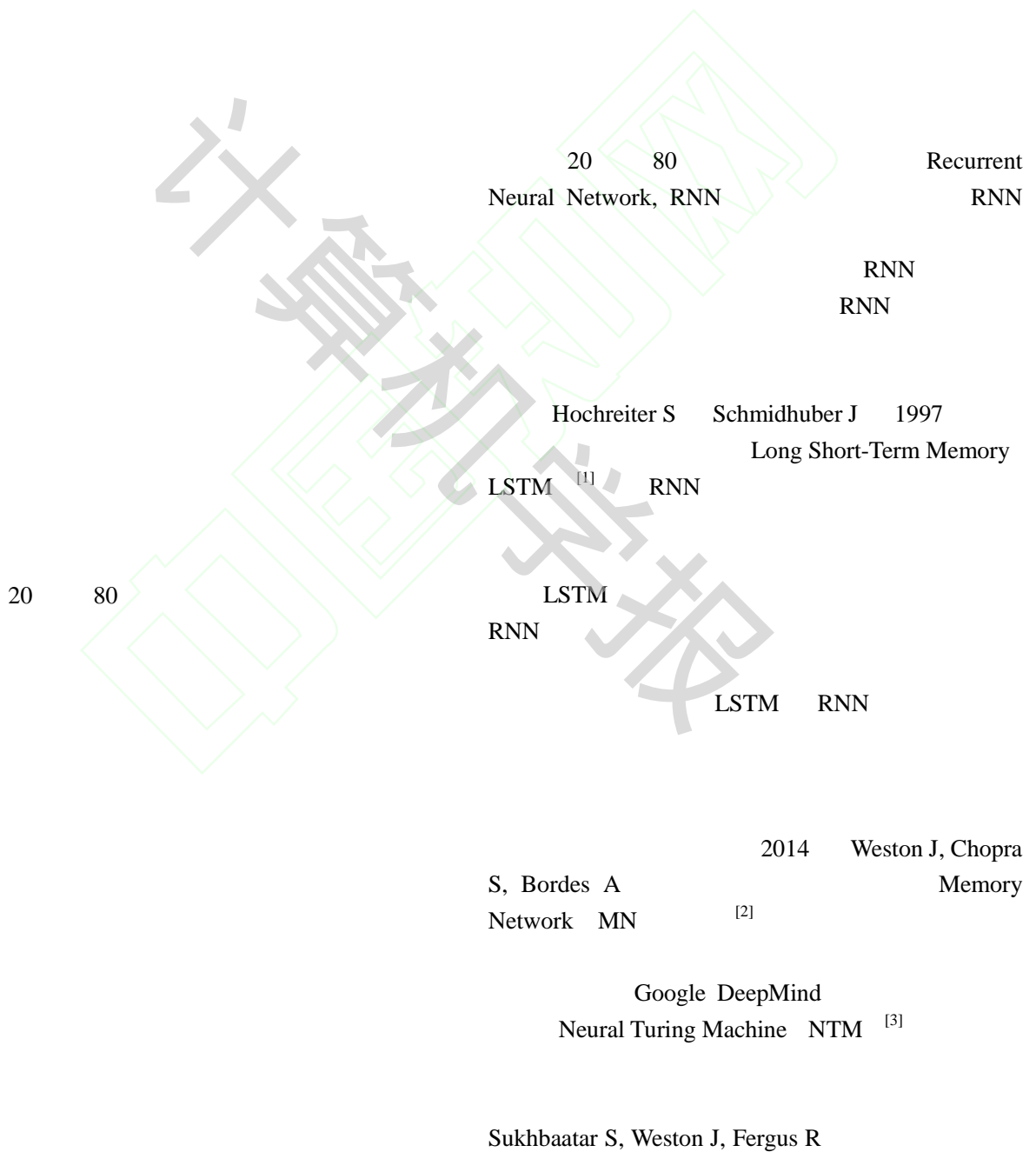
1) 1) 1)

1)

network structure and training algorithm. Afterwards, we introduce the extended model of the memory network and its application in different fields and scenarios. Finally, the future research direction of the memory network is prospected.

Key words recurrent neural network; long short term memory network; memory network; neural turing machine; natural language processing

1



[4]

$$f(Ux_t, Ws_{t-1})$$

2017

Google

Transformer
RNN

[5]

tanh

ReLU

s_{t-1}

$$o_t = \text{soft max}(Vs_t)$$

2

o_t t

2

o_t

t

2014

V

V

RNN

o_t t

[2]

[5][4]

2.1

1

RNN

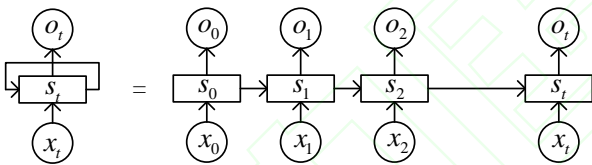
s_t

s_t

t

t

s_t



1 RNN

RNN

U, V, W

RNN

x_t t

x_1

RNN

s_t t

RNN

s_t

2.2

$$s_t = f(Ux_t, Ws_{t-1})$$

1

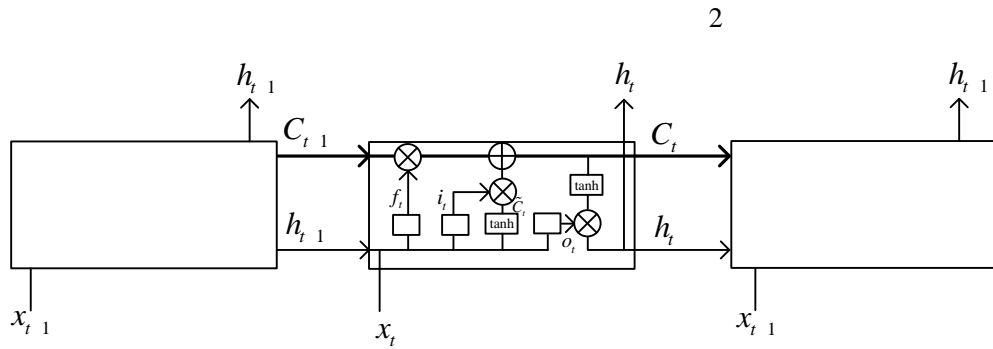
LSTM

RNN

RNN

LSTM RNN

LSTM



2 LSTM

LSTM

2

tanh

LSTM

$$i_t = \text{sigmoid}(W_{i1} [h_{t-1}, x_t] + b_i) \tag{4}$$

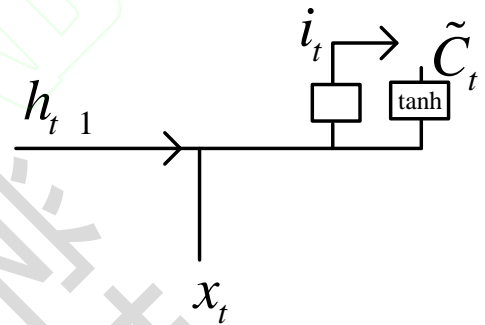
$$\tilde{C}_t = \tanh(W_c [h_{t-1}, x_t] + b_c) \tag{5}$$

Sigmoid



1

LSTM



LSTM

sigmoid

3

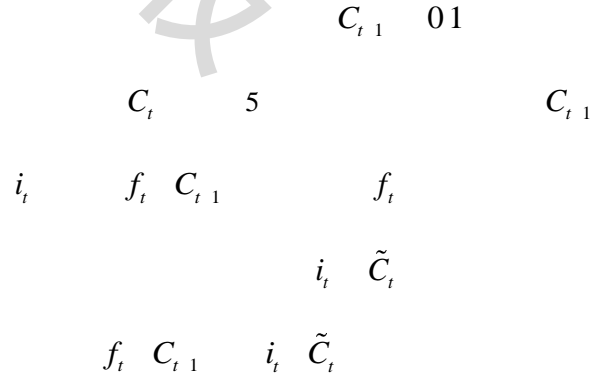
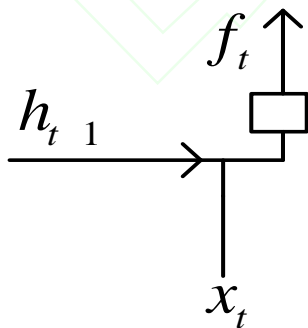
h_{t-1} x_t

f_t

$$f_t = \text{sigmoid}(W_{f1} [h_{t-1}, x_t] + b_f) \tag{3}$$

3

LSTM

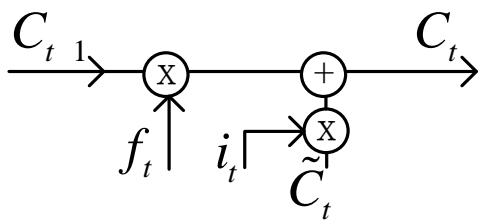


3 LSTM

4

sigmoid

$$C_t = f_t C_{t-1} + i_t \tilde{C}_t \tag{6}$$



$$h_t \quad o_t \quad \tanh(c_t)$$

11

Greff K

[6]

LSTM

LSTM

LSTM

5 LSTM

2.3

C_t

6

sigmoid

2014

Google DeepMind

[3]

o_t

C_t

1936

tanh

C_t

-1 1

NTM

7

o_t

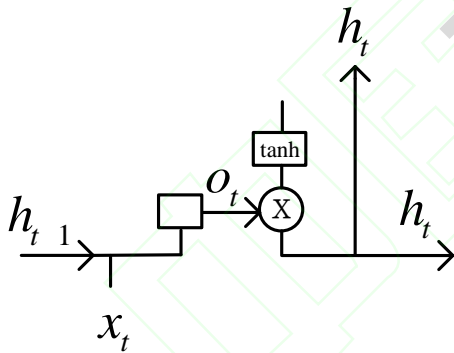
h_t

$$o_t \quad (W_o [h_{t-1}, x_t] + b_o)$$

7

$$h_t \quad o_t \quad \tanh(C_t)$$

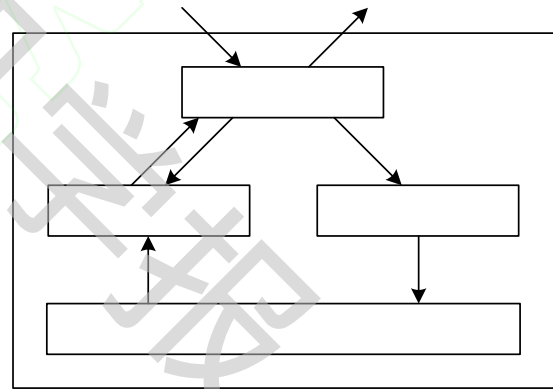
8



6 LSTM

LSTM

c_t



7

NTM

NTM

h_t

t

1

i_t

$$M_t \hat{I}_i^{N \times M} \quad t$$

f_t

$$W \quad h_{t-1}, x_t$$

9

N

M

o_t

$$w_t \hat{I}_i^N \quad t$$

read head

\hat{c}_t

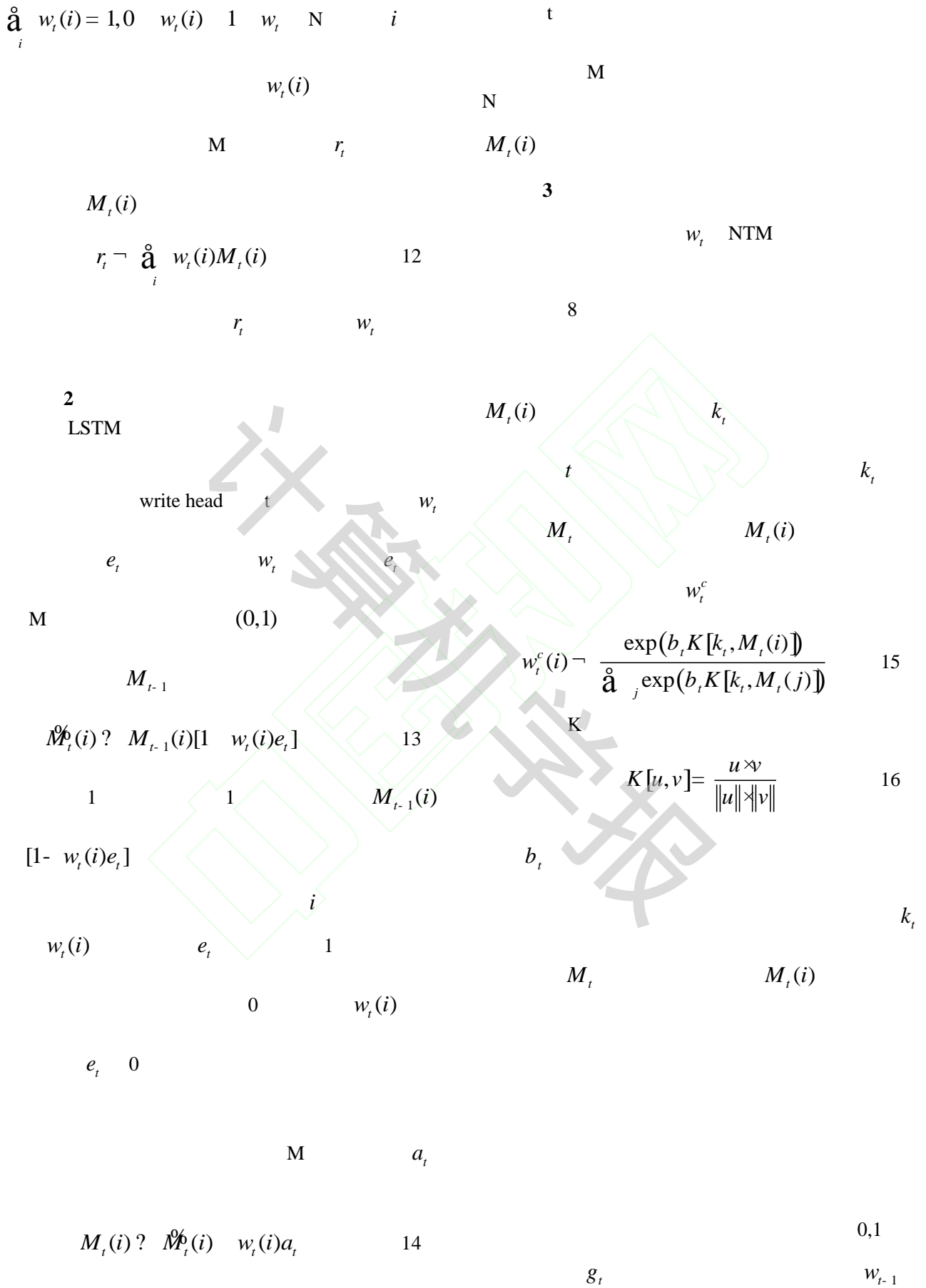
tanh

N

w_t

$$c_t \quad f_t \quad c_{t-1} \quad i_t \quad \hat{c}_t$$

10

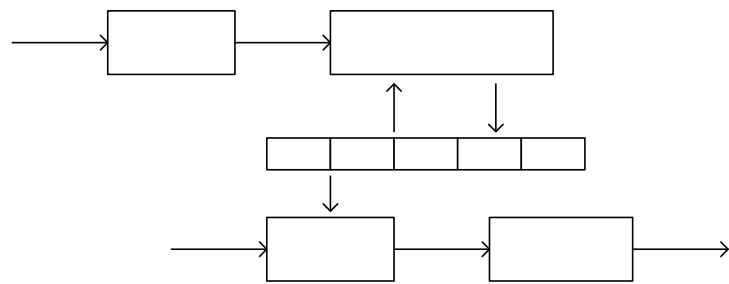
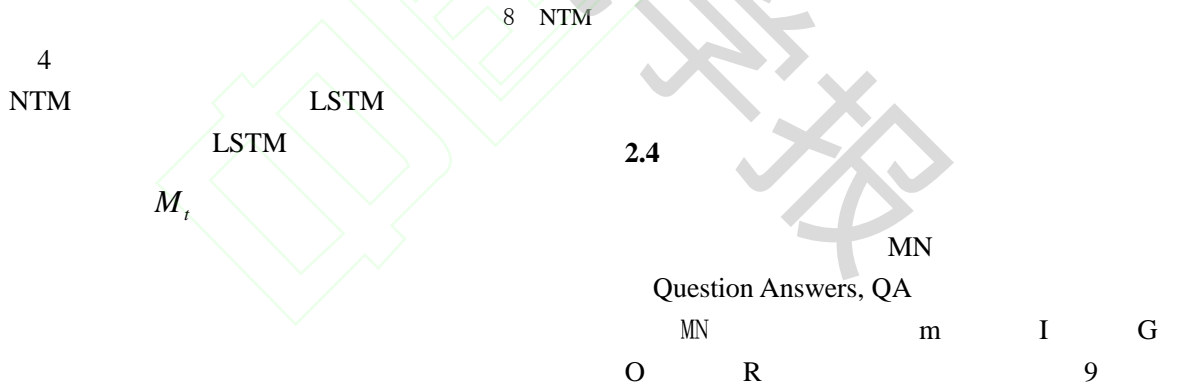
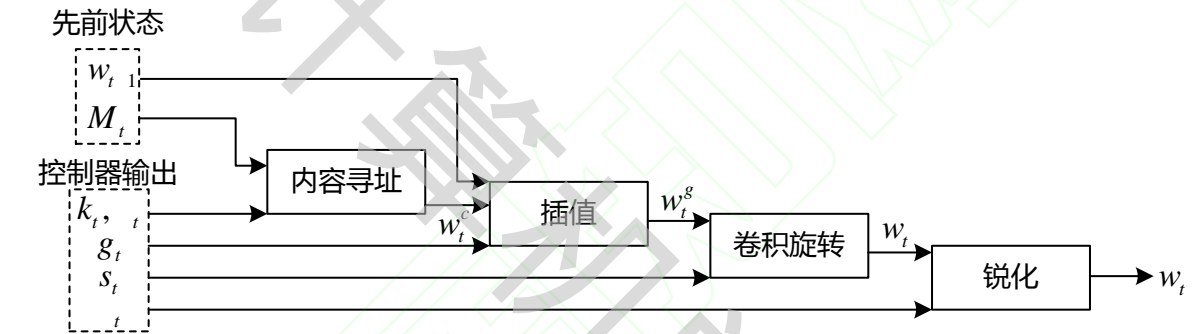


$$w_t^c = g_t w_t^c + (1 - g_t) w_{t-1}^c$$

sharpening

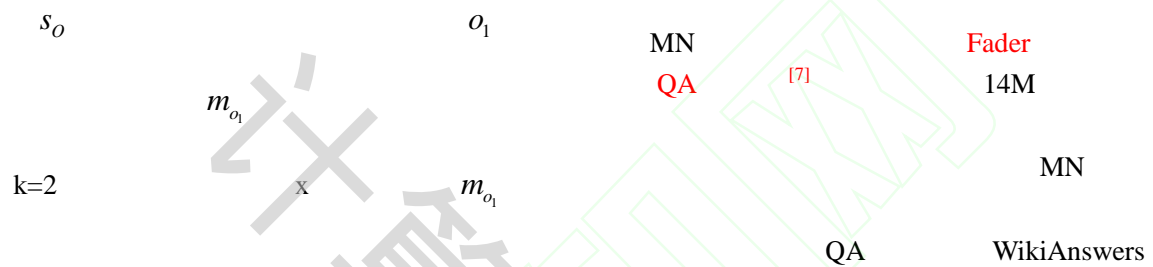
$$w_t(i) = \frac{\phi_t(i)^{g_t}}{\sum_j \phi_t(j)^{g_t}}$$

$$\phi_t(i) = \sum_{j=0}^{N-1} w_t^g(j) s_t(i-j)$$

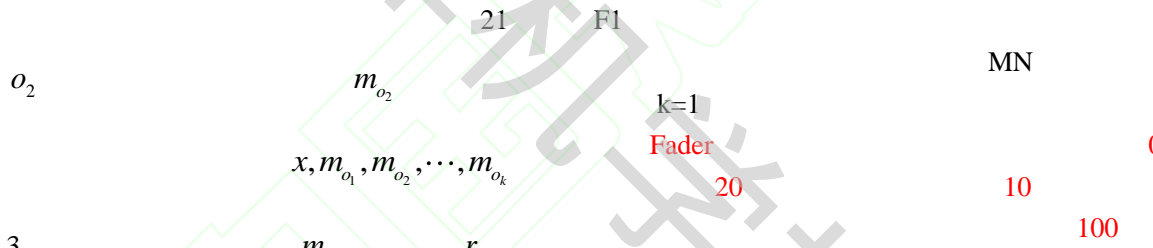


$$\begin{aligned}
 & \max(0, s_O(x, m_{o_1}) - s_O(x, \bar{f})) \\
 I(x) & \max(0, s_O([x, m_{o_1}], m_{o_2}) - s_O([x, m_{o_1}], \bar{f})) \\
 & \max(0, s_R([x, m_{o_1}, m_{o_2}], r) - s_R([x, m_{o_1}, m_{o_2}], \bar{r}))
 \end{aligned}$$

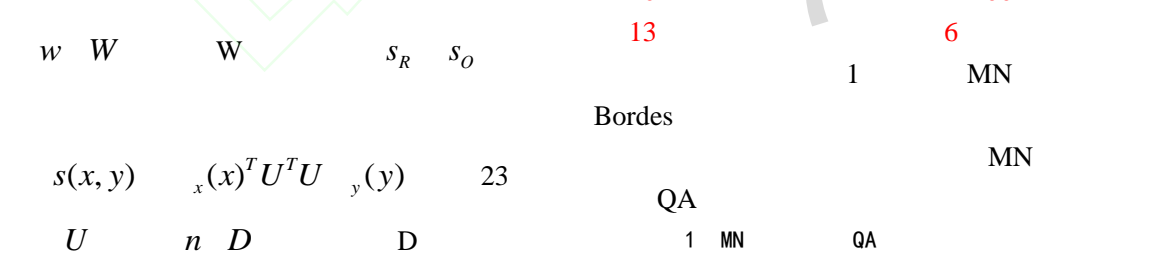
$$o_1 \quad O_1(x, m) \quad \arg \max_{i=1, \dots, N} s_O(x, m_i) \quad 20$$



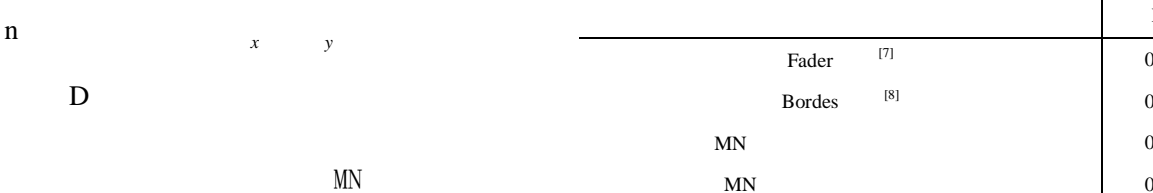
$$o_2 \quad O_2(x, m) \quad \arg \max_{i=1, \dots, N} s_O([x, m_{o_1}], m_i)$$



$$r \quad \arg \max_{w \in W} s_R([x, m_{o_1}, m_{o_2}], w) \quad 22$$



$$s(x, y) = x(x)^T U^T U_y(y) \quad 23$$



margin ranking loss

	1 MN	QA	F1
			0.54
			0.73
	MN		0.72
	MN		0.82

QA :

1 5

		RNN	LSTM
	RNN	LSTM	
		[9]	
	MN		
	100	0.01	0.1
10		2	
2 MN	QA		$\frac{1}{\sqrt{U_2}} \frac{1}{\sqrt{ON_w}}$

$$\begin{matrix}
 & & & & q \\
 k & & v & & \\
 & & \sqrt{d_k} & & d_k \\
 & & & &
 \end{matrix}
 \begin{matrix}
 \sin(pos + k) = \sin(pos)\cos(k) + \cos(pos)\sin(k) \\
 \cos(pos + k) = \cos(pos)\cos(k) - \sin(pos)\sin(k)
 \end{matrix}$$

28

softmax

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

25

WMT 2014 - WMT 2014 - Adam [10]
 $b_1=0.9$ $b_2=0.98$ $e=10^{-9}$

$$\begin{matrix}
 q & & k & & v \\
 & & & & \\
 & & & & \\
 q & & & & k \\
 v & & & & h
 \end{matrix}$$

$$l_{rate} = d_{model}^{-0.5} \min(step_num^{-0.5}, step_num \times warmup_steps^{-1.5})$$

29

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, K, \text{head}_h) W^O$$

Dropout [11]
 $P_{drop} = 0.1$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

26

$$W^O, W_i^Q, W_i^K, W_i^V$$

Position Embedding

RNN

ByteNet

WMT 2014 -

1.2 Dilated Convolution

$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i/d_{model}})$$

27

$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i/d_{model}})$$

pos

i

Deep-Att+PosUnk

WMT 2014 -

256

LSTM

512

d_{model}

GNMT+RL

WMT 2014 -

WMT 2014 -

$pos+k$

pos

8

8

1024

LSTM

			1024	2		
				MoE	MoE	2048
ConvS2S	WMT 2014	-			200	
WMT 2014	-					BLEU
	512					3
	512					
		0.99				
0.1			0.25			
MoE	WMT 2014	-	WMT			4
2014	-					
GNMT						
	LSTM		3			
			3			

	BLEU			
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet ^[13]	23.75			
Deep-Att + PosUnk ^[14]		39.2		1.0 10 ²⁰
GNMT + RL ^[15]	24.6	39.92	2.3 10 ¹⁹	1.4 10 ²⁰
ConvS2S ^[16]	25.16	40.46	9.6 10 ¹⁸	1.5 10 ²⁰
MoE ^[17]	26.03	40.56	2.0 10 ¹⁹	1.2 10 ²⁰
Deep-Att + PosUnk Ensemble ^[14]		40.4		8.0 10 ²⁰
GNMT + RL Ensemble ^[15]	26.30	41.16	1.8 10 ²⁰	1.1 10 ²¹
ConvS2S Ensemble ^[16]	26.36	41.29	7.7 10 ¹⁹	1.2 10 ²¹
Transformer (base model)	27.3	38.1		3.3 10¹⁸
Transformer (big)	28.4	41.0		2.3 10 ¹⁹

4

DTN[18]	2017		
WTN[19]	2017		
AAN[20]	2018		
BlendCNN[21]	2018	CNN	
Action			
Transformer[22]	2018		
Universal			
Transformers[23]	2018		
Evolved			
Transformer[24]	2019		
Set Transformer[25]	2019		
Transformer-XL[26]	2019		

2.6

5

RNN

MN

RNN LSTM NTM

5

1 RNN

MN

2 LSTM

RNN

GPU

MN

3 NTM

RNN

4 MN

NTM

NTM

MN

5

RNN 1986

LSTM 1997

NTM 2014

MN 2014

Transformer 2017

GPU

3

RNN

Attention-based Memory Selection Recurrent
Network AMSRN [49]

RNN LSTM NTM MN

RNN

AMSRN

LSTM

LSTM

MN

LSTM

3.1

RNN

RNN

1 LSTM

LSTM

RNN

RNN

$\{x_1, x_2, \dots, x_t, \dots\}$

x_t

N

1-of-N

RNN

RNN

LSTM

d

3.1.1

RNN

RNN

$h_t R^d$

t

LSTM

$$M_t = [h_0, h_1, \dots, h_{t-1}]$$

2

$$h_t \quad d \quad w_{h1} \quad w_{h2} \quad \text{LSTM}$$

$$M_t = [h_0, h_1, \dots, h_{t-1}] \quad h_t$$

$$d \quad k_t$$

$$k_t \quad W_{kh} h_t \quad b_k \quad 30$$

$$W_{kh} \in \mathbb{R}^{d \times d} \quad b_k \in \mathbb{R}^d$$

$$k_t \quad \text{LSTM}$$

$$M_t = [h_0, h_1, \dots, h_{t-1}] \quad h_t \quad e_{ii}$$

$$e_{ii} = (h_t \circ w_{h1}) \cdot k_t \quad 31$$

o

$$h_t \circ w_{h1} \quad \text{LSTM} \quad M_t = [h_0, h_1, \dots, h_{t-1}]$$

$$h_t \quad w_{h1}$$

$$h_t \circ w_{h1} \quad k_t$$

$$\text{softmax} \quad e_{ii}$$

ii

$$ii \quad \frac{\exp(e_{ii})}{\sum_{i=0}^{t-1} \exp(e_{ii})} \quad 32$$

$$w_{h2} \quad h_t \quad h_t$$

72

$$h_t \quad r_t$$

$$r_t \quad ii \quad h_t$$

$$h_t \quad h_t \circ w_{h2} \quad 33$$

$$r_t = \sum_{i=0}^{t-1} h_{ii}$$

34

$$r_t \quad h_t$$

$$P_w = \text{softmax}(W_{ph} h_t \quad W_{pr} r_t \quad b_p) \quad 35$$

$$W_{ph}, W_{pr}, b_p$$

LSTM+	134.09	93.74	102.04
LSTM+ +	133.36	92.49	86.85
RMN	123.32	64.41	121.28
RMR	134.30	71.04	145.24

7

LSTM
 LSTM
 Penn Treebank
 Switchboard
 Gigaword

$$s_t = (1 - \beta) Bx_t + \beta s_{t-1}$$

$$h_t = (Ps_t, Ax_t, Rh_{t-1})$$

$$y_t = f(Uh_t, Vs_t)$$

A R B P U V

0 1

0 1

Q

t k

$$s_k, k = 1, \dots, t, K$$

Memory Network
 Recurrent [30] AMSRN
 RMN RMR
 RMN RMR
 AMSRN

Penn Treebank Corpus Text8

LSTM RNN

3.1.3

LSTM

RNN

3.1.2

RNN

Mikolov T [31]

RNN

[32]

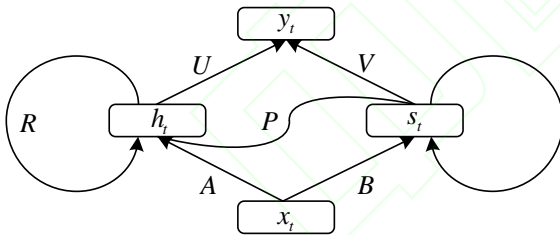
LSTM

1

11

L

12



$$L(D; \theta)$$

$$D = x_i, y_{i-1}^t$$

$$p(D)$$

$$\theta^* = \arg \min_{\theta} E_{D \sim p(D)} [L(D; \theta)]$$

t

y_{t-1}

x_t

h_{t-1}

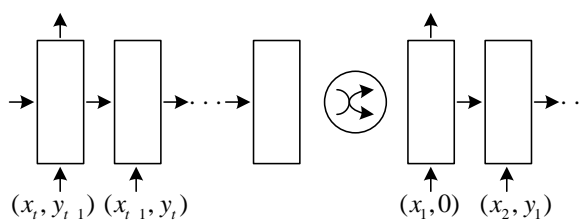
x_t

$$h_t = (Ax_t, Rh_{t-1})$$

36

RNN

s_t



12
2

sigmoid w_{t-1}^r

$$w_{t-1}^{lu} \quad w_t^w \quad () w_{t-1}^r \quad (1 \quad ()) w_{t-1}^{lu} \quad 42$$

$$w_{t-1}^{lu} \quad w_t^w \quad 1 \quad 0$$

k_t

$$M_t(i) \quad M_{t-1}(i) \quad w_t^w(i)k_t \quad 43$$

LSTM
Recently Used Access
LRUA

Least

LRUA

RNN

RNN

Wang J

RNN [33]

w_{t-1}^{lu}

w_t^{lu}

w_t^{lu}

w_{t-1}^{lu}

w_t^r

w_t^r

w_t^w

w_t^r

39

k_t

M_t

$$K(k_t, M_t(i)) = \frac{k_t \cdot M_t(i)}{\|k_t\| \|M_t(i)\|} \quad 40$$

$$w_t^r(i) = \frac{\exp(K(k_t, M_t(i)))}{\sum_j \exp(K(k_t, M_t(i)))} \quad 41$$

$g \quad \|g\|$

$\|g\| > \nu$

[34]

RNN

$$g^{-1} \frac{gv}{\|g\|}$$

44

N_h LSTM

$O(N_h^2)$ LSTM

Danihelka I
 Associative Long Short-Term Memory
 ALSTM ^[36] LSTM Holographic
 Reduced Representations HRR -

HRR -

$$r = (a_r[1]e^{i r[1]}, a_r[2]e^{i r[2]}, \dots)$$

$$\tilde{N}_{h^{(t)}} L \quad y = r \cdot x = (a_r[1]a_x[1]e^{i(r[1] \cdot x[1])}, a_r[2]a_x[2]e^{i(r[2] \cdot x[2])}, \dots)$$

Pascanu R ^[35]

$$r_1, r_2, r_3$$

46

$$W = \mathring{a}_t \left\| \frac{(\tilde{N}_{h^{(t)}} L) \left\| \frac{h^{(t)}}{h^{(t-1)}} \right\|}{\tilde{N}_{h^{(t)}} L} \right\|$$

45

RNN

3.2 LSTM
 LSTM
 LSTM
 LSTM LSTM LSTM
 LSTM LSTM

$$h = \begin{bmatrix} h_{real} \\ h_{imaginary} \end{bmatrix} \quad 48$$

$$r_{i,s} = \begin{bmatrix} P_s & 0 \\ 0 & P_s \end{bmatrix} r_i \quad 54$$

$$h \in \mathbb{R}^{N_h}, h_{real}, h_{imaginary} \in \mathbb{R}^{N_h/2} \quad c_{s,t} = g_f \circ c_{s,t-1} + r_{i,s} (g_i \circ u) \quad 55$$

LSTM

$$\hat{r}_i, \hat{r}_o$$

$$r_{i,s} = \begin{bmatrix} P_s & 0 \\ 0 & P_s \end{bmatrix} \in \mathbb{R}^{N_h/2 \times N_h/2}$$

s

$$\hat{g}_f, \hat{g}_i, \hat{g}_o, \hat{r}_i, \hat{r}_o = W_{xh} x_t + W_{hh} h_{t-1} + b_h \quad 49$$

$$r = u \begin{bmatrix} r_{real} \circ u_{real} & r_{imaginary} \circ u_{imaginary} \\ r_{real} \circ u_{imaginary} & r_{imaginary} \circ u_{real} \end{bmatrix} \quad 56$$

$$\hat{u} = W_{xu} x_t + W_{hu} h_{t-1} + b_u \quad 50$$

$$r_{o,s} = r_{i,s}$$

0 1

$$\text{bound}(h) = \begin{bmatrix} h_{real} / d \\ h_{imaginary} / d \end{bmatrix} \quad 51$$

$$r_{o,s} = \begin{bmatrix} P_s & 0 \\ 0 & P_s \end{bmatrix} r_o$$

d

$$\max(1, \sqrt{h_{real} \circ h_{real} + h_{imaginary} \circ h_{imaginary}})$$

$\in \mathbb{R}^{N_h/2}$

52

$$h_t = g_o \circ \text{bound}\left(\frac{1}{N_{copies}} \sum_{s=1}^{N_{copies}} r_{o,s} c_{s,t}\right) \quad 58$$

d /

3.2.2 LSTM

Zhang X

[37]

Tree Long Short-Term Memory Networks

TLSTM LSTM

r_i

r_o

$$u = \text{bound}(\hat{u})$$

$$r_i = \text{bound}(r_i)$$

$$r_o = \text{bound}(r_o)$$

53

TLSTM

$D(w)$

$$r_i \in \mathbb{R}^{N_h}$$

w

$$r_o \in \mathbb{R}^{N_h}$$

$$w_0$$

1', 2', ..., n

u

$$g_i$$

LEFT

$$w_0$$

$$w_k$$

$$w_0$$

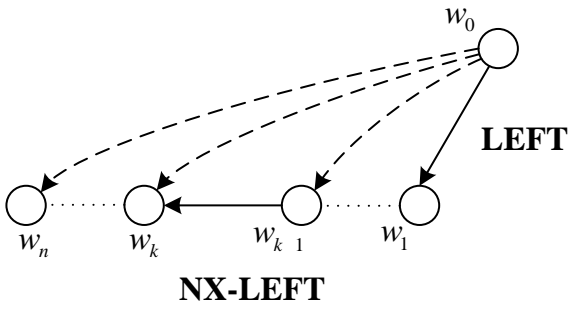
$$w_{k-1} \quad w_k$$

s = {1, ..., N_copies}

NX-LEFT

13

RIGHT NX-RIGHT



13

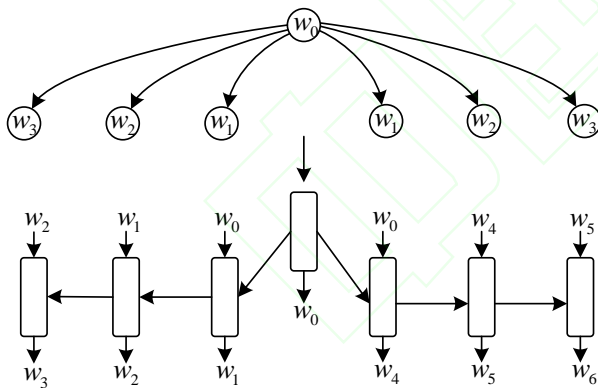
$t \in \{1, n\}$ $\langle w_t, z_t \rangle$ $D(w_t)$
 $t \in \{1, \dots, n\}$ w
 $z_t \in \{\text{LEFT, RIGHT, NX-LEFT, NX-RIGHT}\}$
 LSTM
 $W \in \mathbb{R}^{d \times d}$ $H \in \mathbb{R}^{d \times n}$
 $t \in \{1, \dots, n\}$ $H[:, t]$

T
 S
 breadth-first search BFS
 ROOT

14

T S w
 $P(S|T)$ $P(w|D(w))$ 59
 $w \in \text{BFS}(T) \setminus \text{ROOT}$

$D(w)$
 $x_t \in \mathbb{R}^d$ $W_e \in \mathbb{R}^{d \times d}$ $e(w_t) \in \mathbb{R}^d$
 $h_t \in \mathbb{R}^d$ $\text{LSTM}^{z_t}(x_t, H[:, t])$
 $H[:, t] \in \mathbb{R}^{d \times t}$ $h_t \in \mathbb{R}^d$
 $y_t \in \mathbb{R}^d$ $W_{ho} \in \mathbb{R}^{d \times d}$ $h_t \in \mathbb{R}^d$
 $W_e \in \mathbb{R}^{d \times |V|}$
 $W_{ho} \in \mathbb{R}^{d \times |V|}$
 $H \in \mathbb{R}^{d \times (n-1)}$
 $s \in \{1, \dots, |V|\}$
 $d \in \{1, \dots, |V|\}$



14 TLSTM

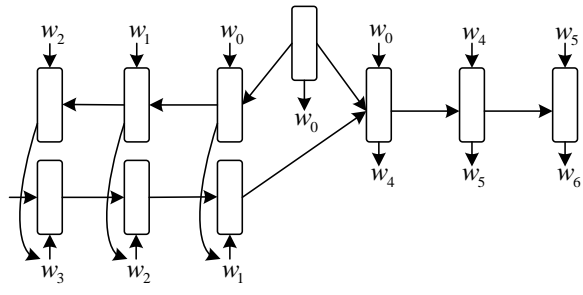
$e(w_t) \in \mathbb{R}^d$ $w_t \in \mathbb{R}^d$ LSTM
 $h_t \in \mathbb{R}^d$ $x_t \in \mathbb{R}^d$ LSTM
 $D(w_t) \in \mathbb{R}^d$ softmax

$$P(w_t | D(w_t)) = \frac{\exp(y_t, w_t)}{\sum_{k=1}^{|V|} \exp(y_t, k)}$$

61
 TLSTM

$D(w) <$ LSTM
 $>$ LSTM
 LSTM LEFT LSTM
 NX-LEFT RIGHT NX-RIGHT
 LSTM

15



15

w_0

中国知网

$$p_j^t = \text{soft max}(u^T \tanh(W_j x_t - W_{j-1} \tilde{x}_{t-1}))$$

72

$$P(\hat{y} | \hat{x}, S) = P(\hat{y} | \hat{x}, S)$$

1

$$\tilde{x}_t = \sum_{j=1}^m p_j^t x_j$$

73

$$\hat{y} = \prod_{i=1}^k a(\hat{x}, x_i) y_i$$

r_t

$$P(\hat{y} | \hat{x}, S)$$

$$r_t = (W_r [\tilde{x}_t, x_t])$$

74

$$c_s(\hat{x}) = x_i, y_i$$

a

\hat{x}

c_t

x_i

a

h_t

c softmax

$$c_t = r_t \tilde{x}_t f_t \tilde{c}_t i_t \hat{c}_t$$

75

$$a(\hat{x}, x_i) = \frac{e^{c(f(\hat{x}), g(x_i))}}{\sum_{j=1}^k e^{c(f(\hat{x}), g(x_j))}}$$

77

$$h_t = o_t \tanh(c_t)$$

76

Penn Treebank
RNN LSTM

LSTM

c

Stanford Sentiment Treebank
LSTM

2

$$g(x_i, S)$$

LSTM

3.2.4

T-CNN
LSTM

$$LSTM(g(x_i, S), \vec{h}_i, \vec{h}_i, g(x_i))$$

LSTM

S

$$g(x_i)$$

$$x_i$$

[40]

$$h_i, c_i = LSTM(g(x_i), h_{i-1}, c_{i-1})$$

78

$$h_i, c_i = LSTM(g(x_i), h_{i-1}, c_{i-1})$$

k

-

h_i

c_i

LSTM

$$S = \{(x_i, y_i)\}_{i=1}^k$$

$$c_s(\hat{x})$$

\vec{h}

$$i = |S|$$

\hat{x}

\hat{x}

\hat{y}

3

\hat{x}

$$f(\hat{x}, S)$$

$$S = c_s(\hat{x})$$

LSTM attLSTM(, ,)

$f(\hat{x}, S)$ attLSTM($f(\hat{x}), g(S), K$) 79

Hierarchical Attentive Memory HAM

f g

LSTM

MN

DMN

K LSTM

KV-MemNN

HMN

$g(S)$

S

x_i

3.3.1

MN

k

[4]

\hat{h}_k, c_k LSTM($f(\hat{x}), [h_{k-1}, r_{k-1}], c_{k-1}$)

h_k \hat{h}_k $f(\hat{x})$

$r_{k-1} = \sum_{i=1}^{|S|} a(h_{k-1}, g(x_i))g(x_i)$

80

x_1, \dots, x_n

q

a

$a(h_{k-1}, g(x_i)) = \text{softmax}(h_{k-1}^T g(x_i))$

LSTM(x, h, c)

x

h

a

19

c

a

1

x_i

$g(S)$

r_{k-1}

h_{k-1}

m_i

c_i

LSTM attLSTM($f(\hat{x}), g(S), K$)

h_K

m_i

u

3.3

MN

MN

MN

p_i

p_i

MN

MN

c_i

o

o

\hat{a}

u

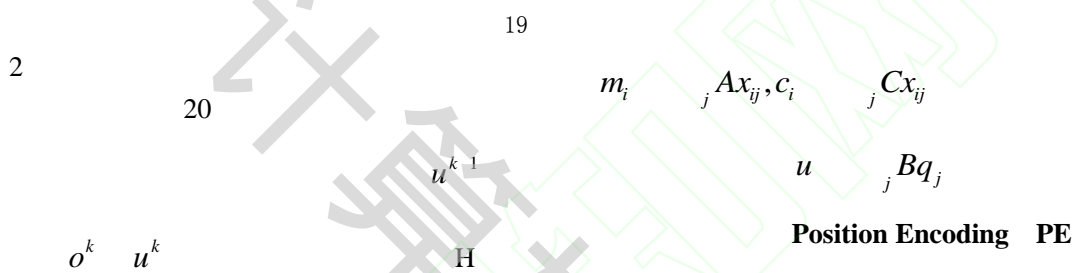
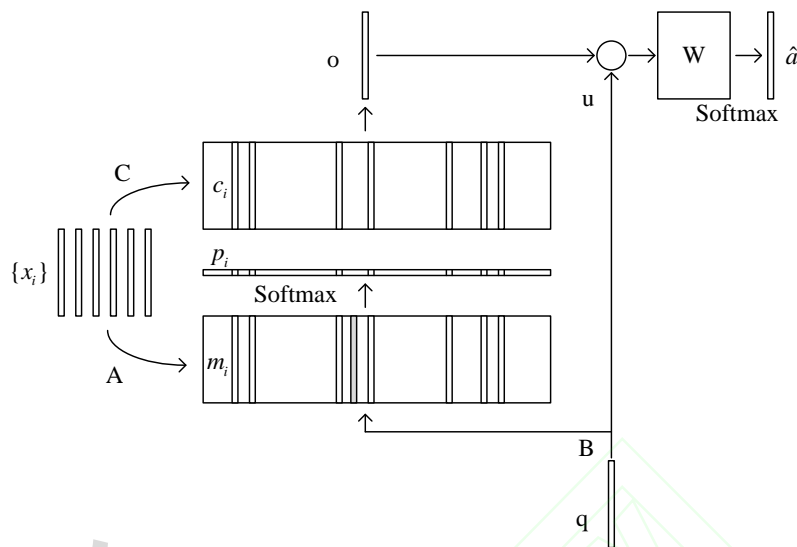
W

Dynamic Memory Networks DMN -

Key-Value Memory Networks KV-MemNN

Hierarchical Memory Network

HMN



$$u^{k-1} \quad H \quad u^k \quad o^k$$

$$u^{k-1} \quad Hu^k \quad o^k$$

$$m_i \quad l_j \quad Ax_{ij}$$

$$l_{kj} \quad (1 - j/J) \quad (k/d)(1 - 2j/J)$$

$\{x_i\}$

W

J

d

$T_A \quad T_C$

\hat{a}

$$\hat{a} = \text{soft max}(Wu^{k-1}) \quad \text{soft max}(W(o^k \quad u^k))$$

82

$$m_i \quad l_j \quad Ax_{ij} \quad T_A(i) \quad 84$$

Bag of Words BoW

$$c_i \quad j \quad Cx_{ij} \quad T_C(i) \quad 85$$

$x_{ij} \quad i \quad j$

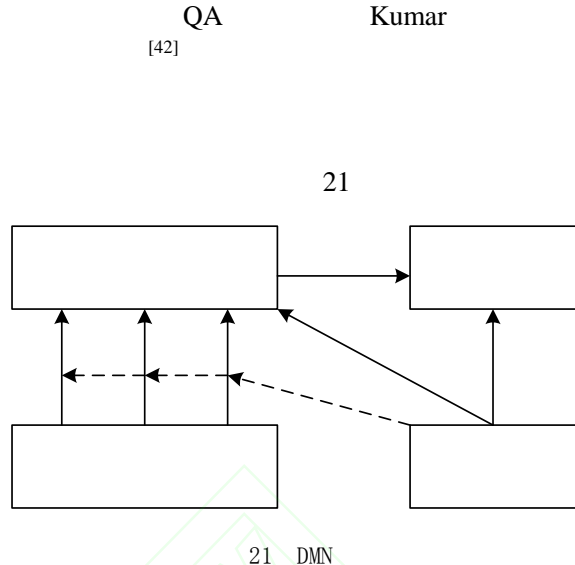
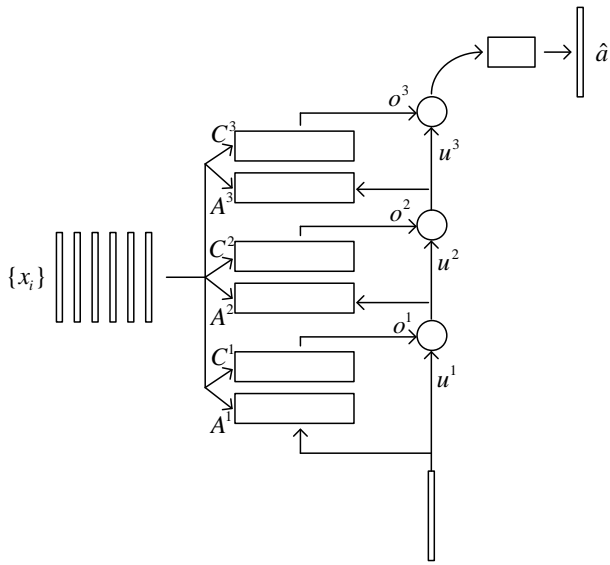
i

,

\hat{a}

$x_i \quad \{x_{i1}, x_{i2}, \dots, x_{in}\}$

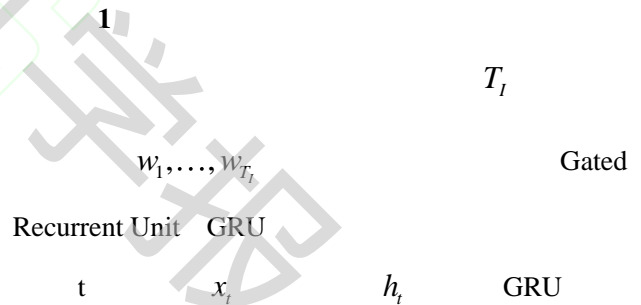
a



20 Penn Treebank Text8 Structurally
 RNN LSTM Constrained Recurrent Nets, SCRN

3.3.2

[41]



u^{k-1}

$$T^k(u^k) = (W_T^k u^k + b_T^k) \quad 86$$

$$u^{k-1} \circ o^k = T^k(u^k) \circ u^k = (1 - T^k(u^k)) \circ u^k \quad 87$$

$$W_T^k, b^k \quad k$$

$$T^k \quad k$$

$$z_t = (W^{(z)} x_t + U^{(z)} h_{t-1} + b^{(z)})$$

$$r_t = (W^{(r)} x_t + U^{(r)} h_{t-1} + b^{(r)})$$

$$\tilde{h}_t = \tanh(W x_t + r_t \circ U h_{t-1} + b^{(h)}) \quad 88$$

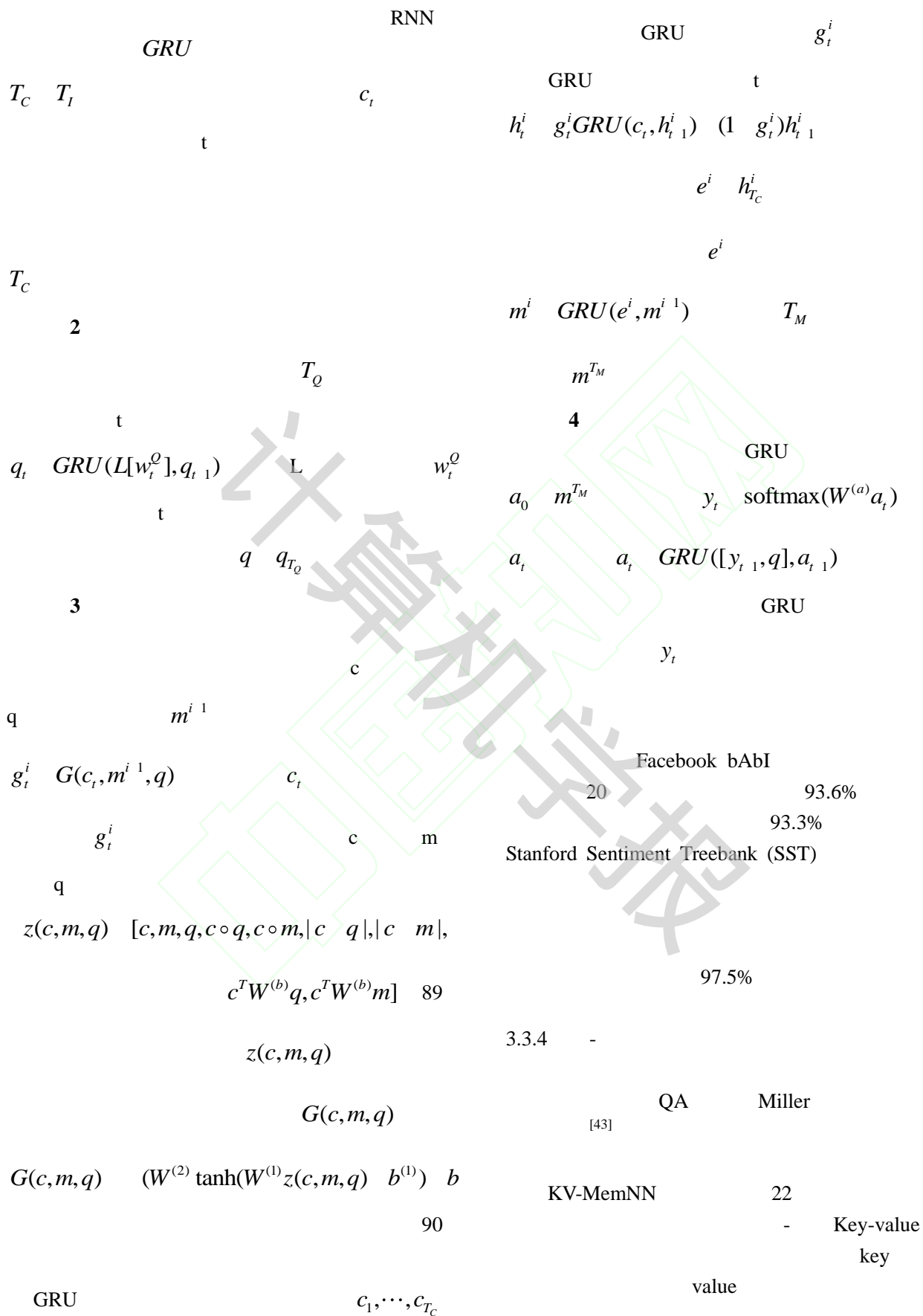
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

$$W^{(z)}, W^{(r)}, W \in \mathbb{R}^{n_H \times n_I} \quad U^{(z)}, U^{(r)}, U \in \mathbb{R}^{n_H \times n_H}$$

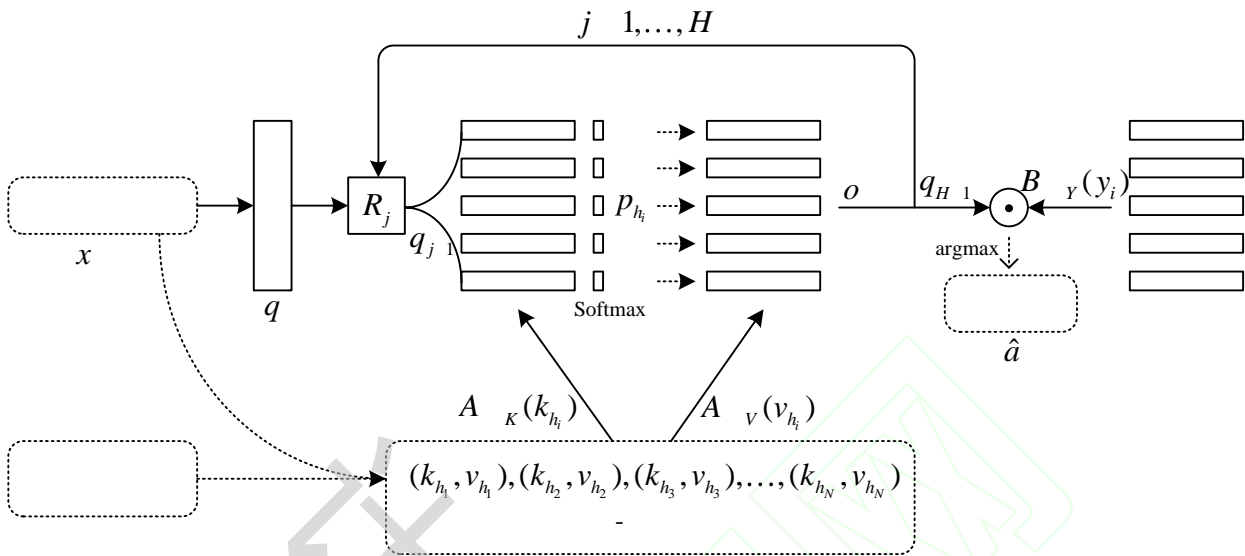
n GRU

3.3.3

$$h_t = GRU(x_t, h_{t-1})$$



$(k_1, v_1), \dots, (k_M, v_M)$



22 KV-MemNN

1 Key Hashing

j hop R_j

N

$(k_{h_1}, v_{h_1}), \dots, (k_{h_N}, v_{h_N})$

1000

q

$p_{h_i} \text{ softmax}(q_{j+1}^T A_K k_{h_i})$

2 Key Addressing

H

Softmax

k_{h_i}

x

\hat{a}

$\text{argmax}_{i=1, \dots, C} \text{softmax}(q_{H-1}^T B_Y(y_i))$

91

$p_{h_i} \text{ Softmax}(A_X(x) A_K(k_{h_i}))$ 93

y_i

Y

D

A

$d' D$

3 Value Reading

3.3.5

$o_i p_{h_i} A_V(v_{h_i})$

Chandar

Maximum Inner Product Search

x

MIPS

Hierarchical Memory Network HMN

[44]

$q A_X(x)$

$o q$

$q_2 R_1(q o)$

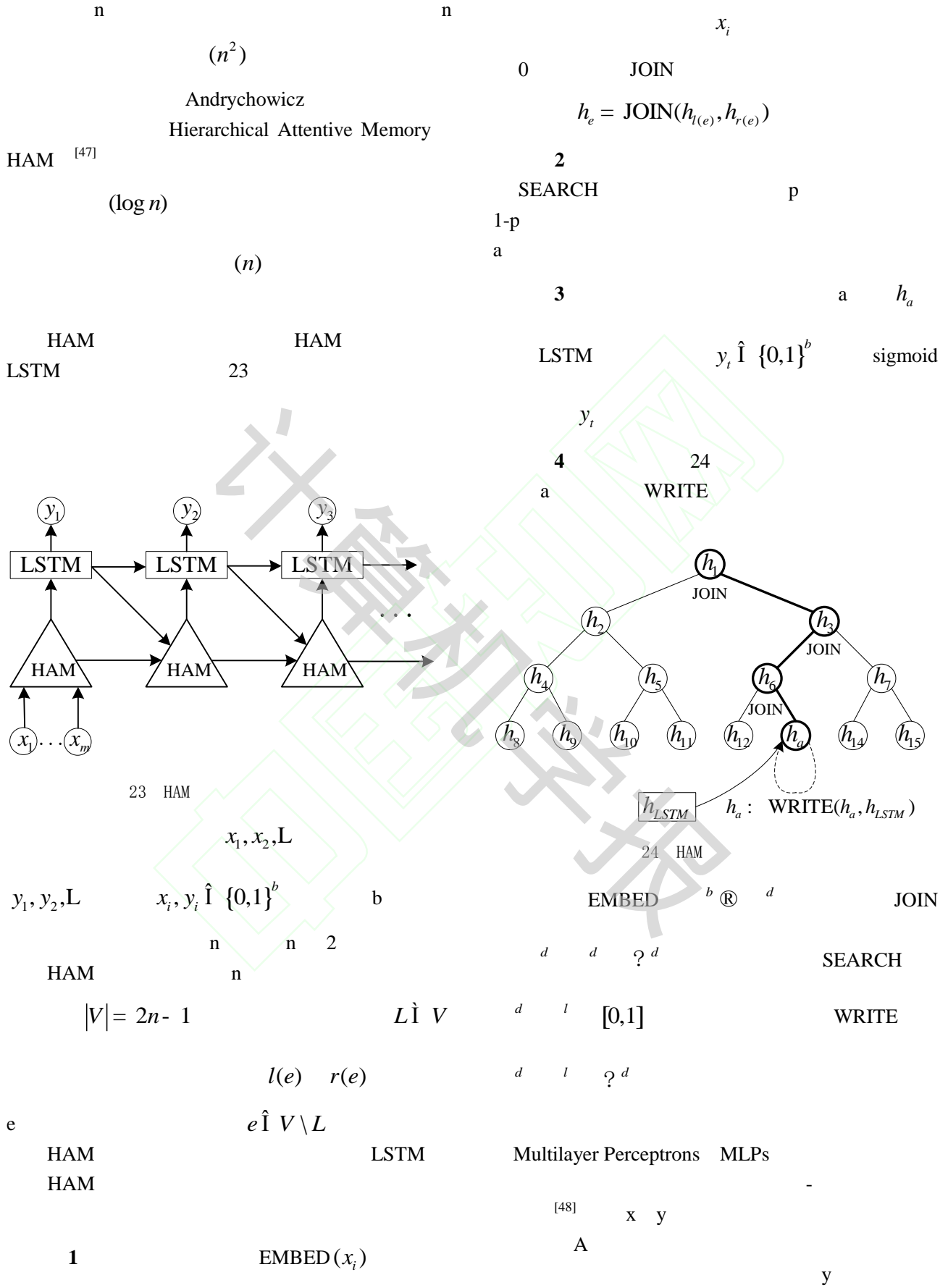
R_1

$d d$

MN

HMN

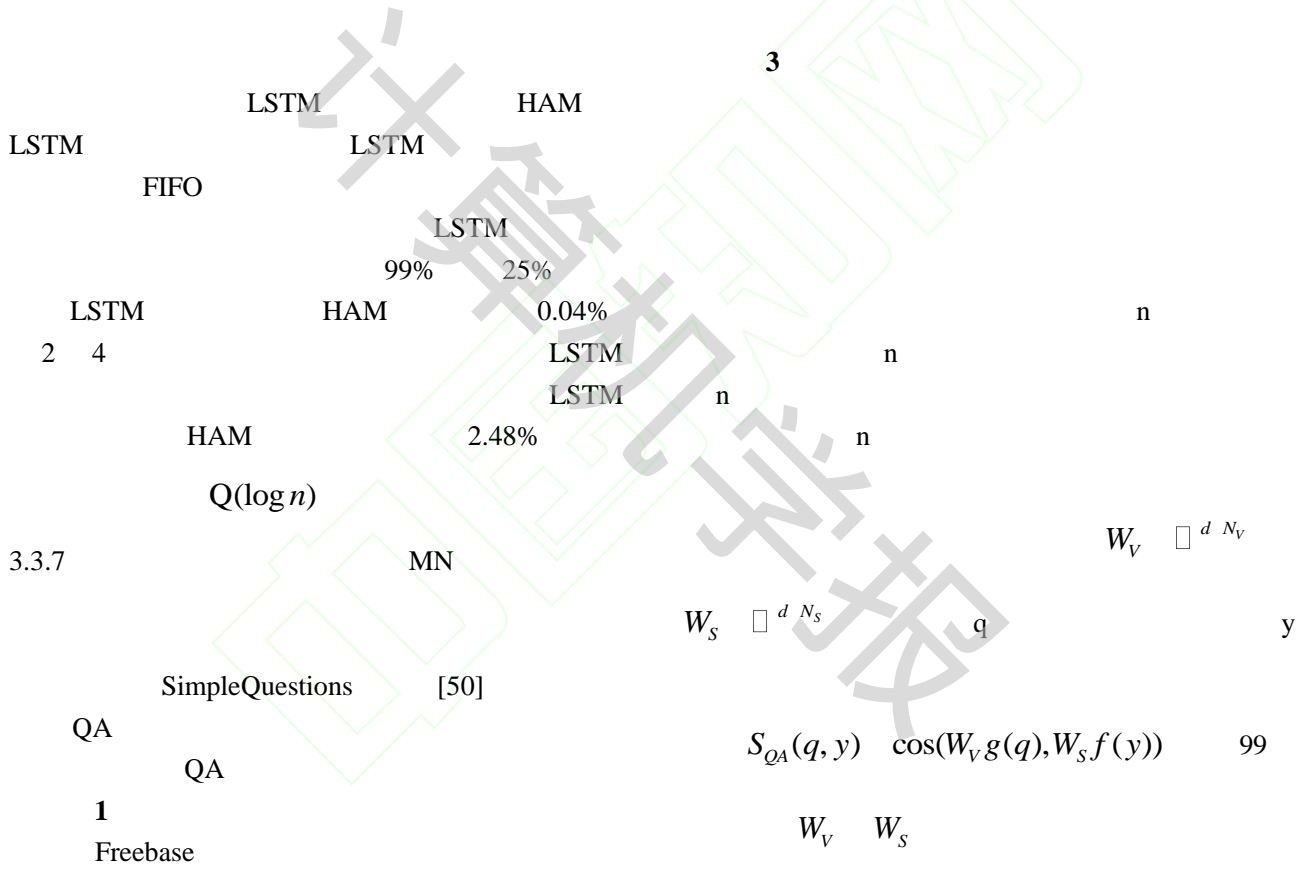
1	HMN			MIPS		Maximum
				Cosine Similarity Search	MCSS	
	hash			$\operatorname{argmax}_i^{(K)} \frac{q^T x_i}{\ q\ \ x_i\ }$	$\operatorname{argmax}_i^{(K)} \frac{q^T x_i}{\ x_i\ }$	93
		MIPS			x_i	MCSS
MIPS				MIPS		
					MIPS	MCSS
2	HMN		MIPS		P Q	MCSS
					MIPS	MCSS
MIPS				$P(x)$	$x, 1/2$	$\ x\ _2^2, 1/2$
				$Q(x)$	$x, 0, 0, \dots, 0$	$\ x\ _2^4, \dots, 1/2$
						$\ x\ _2^{2m}$
						94
HMN	K-MIPS				q	MIPS
$\{x_1, \dots, x_n\}$		q	K-MIPS			
$\operatorname{argmax}_i^{(K)} q^T x_i$		$\operatorname{argmax}^{(K)}$	K	$\operatorname{argmax}_i^{(K)} q^T x_i$	\square	$\operatorname{argmax}_i^{(K)} \frac{Q(q)^T P(x_i)}{\ Q(q)\ _2 \ P(x_i)\ _2}$
$q^T x_i$		$\{x_1, \dots, x_n\}$				95
		q		MIPS	MCSS	K
				[46]		
HMN	K-MIPS			K	K	
softmax	MIPS					HMN
C	$\operatorname{argmax}^{(K)} h(q) M^T$			1	K	
R_{out}	$\operatorname{softmax}^{(K)}(h(q) M^T)$	92		mini-batch		
	$\operatorname{softmax}(h(q) M [C]^T)$			GPU		
$h(q)$	\square^d	C	K	2		
MIP		M	$\square^{N \times d}$			
N		$M[C]$	M	3		
				K		
$M[C]$	C			3.3.6		MN
	K-MIPS		Aurolat			
[45]						
K-MIPS			MIPS			(n)



$$L = \log p(y|x, A) = \log \prod_A p(A|x, A) p(y|A, x, A)$$

$$F_A = \prod_A p(A|x, A) p(y|A, x, A) \quad L \quad 97$$

$$F_A = \prod_A p(A|x, A) [\log p(y|A, x, A) - \log p(y|A, x, A)] \quad 98$$



$$W_S \in \mathbb{R}^{d \times N_S} \quad q$$

$$W_V \in \mathbb{R}^{d \times N_V} \quad y$$

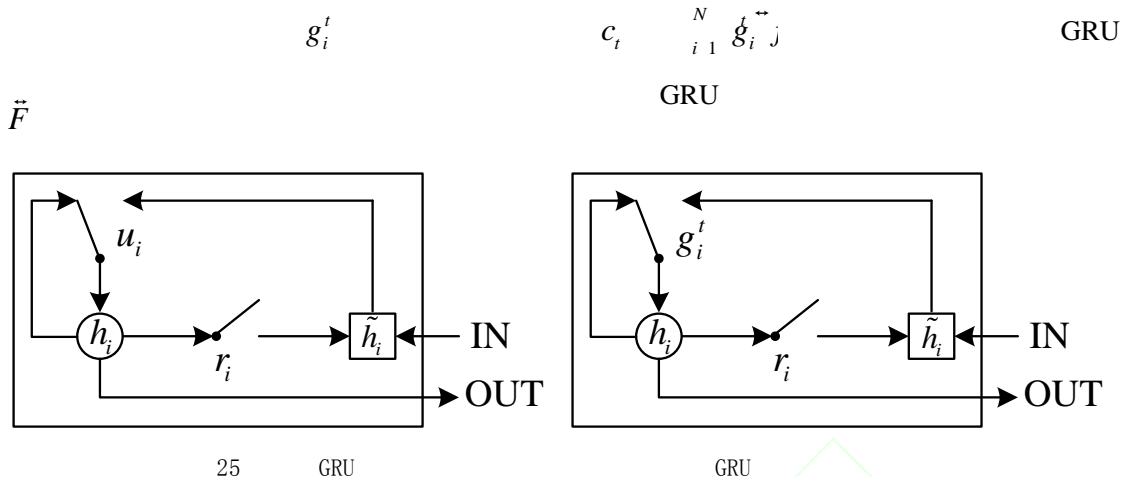
$$S_{QA}(q, y) = \cos(W_V g(q), W_S f(y)) \quad 99$$

$$y = (s, r, \{o_1, \dots, o_k\}) \quad f(y) \in \mathbb{R}^{N_S} \quad S_{RVB}(q, y) = \cos(W_V g(q), W_{VS} h(y)) \quad 100$$

$$N_S \quad f(y) \quad W_V \quad W_S$$

1/k

1



25 GRU

$h_i = g_i^t \circ \tilde{h}_i$ u_i

$\tilde{h}_i = (1 - g_i^t) \circ h_{i-1}$ 105

3.3.9 KV- Rama()-249()-238()JTJEMC 12

h_N c_t

m^t m^{t-1}

$m^t = GRU(c^t, m^{t-1})$ 106

GRU

ReLU

$m^t = RuLU(W^t[m^{t-1}; c^t; q] - b)$ 107

$W^t \in \mathbb{R}^{n_H \times n_H}, b \in \mathbb{R}^{n_H}, n_H$

Ramachandran DMN

Tensor Network DMTN [52]

Dynamic Memory [53]

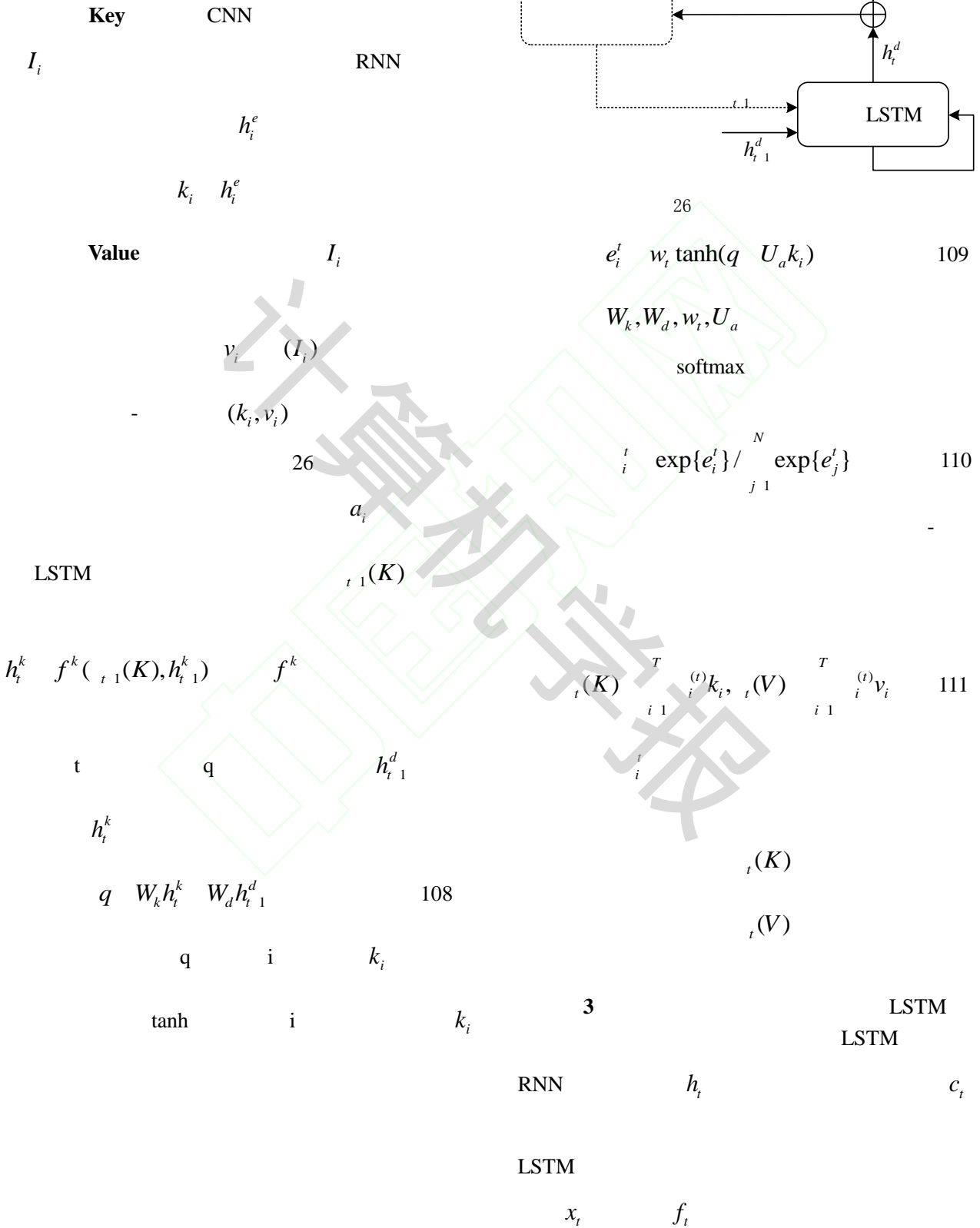
DMN 80%

Facebook

bAbI 20%

Visual Question Answering VQA

$(k_1, v_1), \dots, (k_T, v_T)$



[55]

o_t

h_t

softmax

k

k

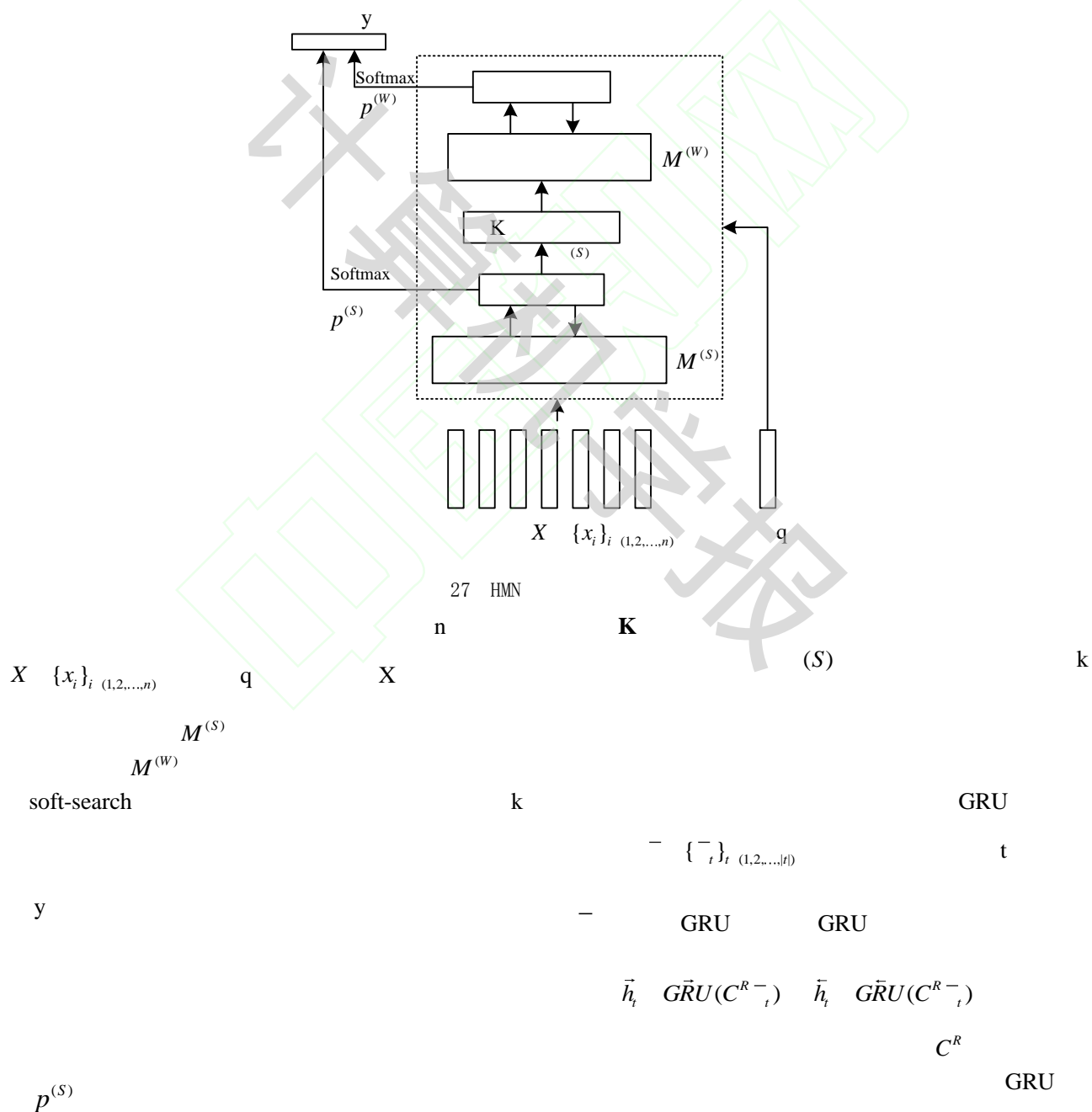
p_t

$$p_t = \text{softmax}(U_p[h_t, x_t, \dots, v_t(V)]) + b_p \quad 112$$

3.3.10

HMN

27



\bar{h}_t GRU \bar{h}_t [56]

$M^{(w)} \{m_t\}_{t=(1,2,\dots,t)}$ $m_t \bar{h}_t \bar{h}_t$

m_t X t

\bar{w}_t

HMN
HMN

$\hat{\{ \hat{m}_t \}}_{t=(1,2,\dots,\hat{t})}$

\hat{X}

3.4

$u_R^{(S)}$

\hat{w}

$\{\hat{m}_t\}_{t=(1,2,\dots,\hat{t})}$

$\{m_t^{(w)}\}_{t=(1,2,\dots,\hat{t})}$

3.4.1

$t^{(w)} \text{softmax}(v^T \tanh(Wu_R^{(S)} U\hat{m}_t))$ 113

$v^{d \times 1}, W^{d \times d}, U^{d \times d}$

urlanello
Active Long Term Memory

Networks A-LTM [58]

Oriol Vinyals [53]

Distillation Loss [59]

\hat{w}

McClelland
Hippocampus

[60]
Neocortex

$p^{(w)}(\cdot) \text{trans}(p^{(w)}(\hat{\cdot})) \text{trans}(m^{(w)})$ 114

$\text{trans}(\cdot)$

Catastrophic Inference

$p^{(w)}(\hat{w}) \hat{I}_i^{|\hat{I}|}$

CI

$p^{(w)}(w) \hat{I}_i^{|\hat{I}|}$

A-LTM

0

N

$M^{(S)}$

H

$p^{(S)}$

$M^{(w)}$

$p^{(w)}$

H
N

N

A-LTM

$p(w) = p^{(S)}(w) + p^{(w)}(w)$ 115

N

y

N

N

0

$$w_0^0 \quad w_0^* \quad w_1^0 \quad w_1^* \quad w_2^0 \sim N(0, \sigma^2) \quad 116$$

$$f(w_0^*, w_1^*; x_1)$$

N

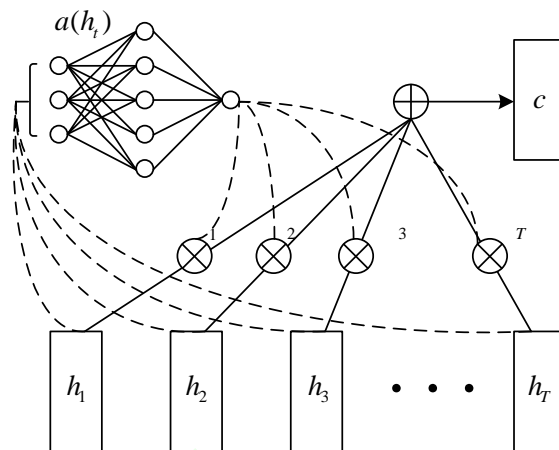
y_1

3.4.2

Colin Raffel

[61]

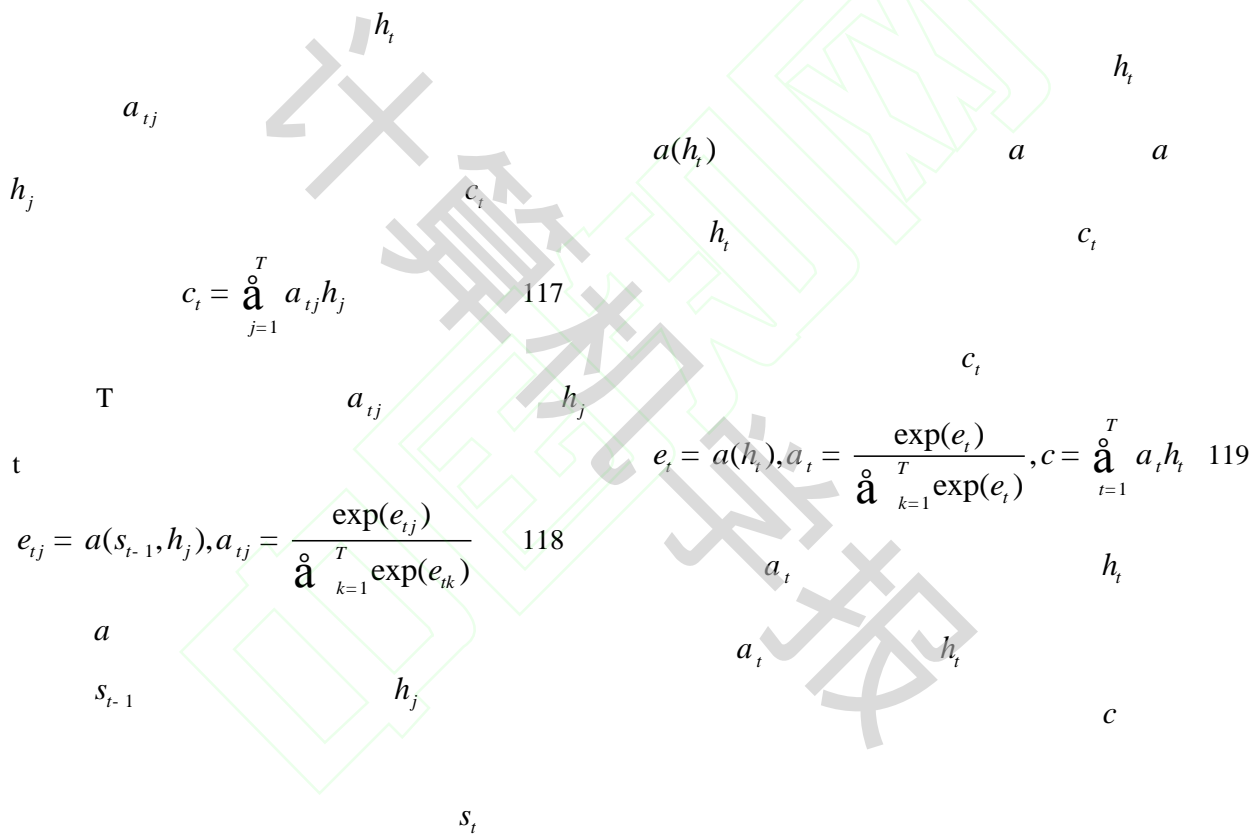
28



Bahdanau

[62]

28



h_j

a_{ij}

h_t

c_i

$a(h_t)$

h_t

h_t

a

h_t

a

c_i

$$c_i = \mathring{\mathbf{a}} \sum_{j=1}^T a_{ij} h_j \quad 117$$

T

a_{ij}

h_j

c_i

t

$$e_t = a(h_t), a_t = \frac{\exp(e_t)}{\mathring{\mathbf{a}} \sum_{k=1}^T \exp(e_k)}, c = \mathring{\mathbf{a}} \sum_{t=1}^T a_t h_t \quad 119$$

$$e_{ij} = a(s_{t-1}, h_j), a_{ij} = \frac{\exp(e_{ij})}{\mathring{\mathbf{a}} \sum_{k=1}^T \exp(e_{ik})} \quad 118$$

a

s_{t-1}

h_j

a_t

a_t

h_t

h_t

c

s_t

s_t

s_{t-1}

T

T

c_t

t-1

h_t

$$s_t = f(s_{t-1}, c_t, y_{t-1})$$

$$c = \frac{1}{T} \mathring{\mathbf{a}} \sum_{t=1}^T h_t$$

x_t

$$h_t = \text{LReLU}(W_{xh}x_t + b_{xh}) \quad 120$$

$$y = \text{LReLU}(W_{sy}s + b_{sy}), W_{sy} \in \mathbb{R}^{1 \times D}, b_{sy} \in \mathbb{R}^1$$

$$W_{xh} \in \mathbb{R}^{D \times 2}, b_{xh} \in \mathbb{R}^D, \text{LReLU}(x)$$

122

y
adam

$$\text{LReLU}(x) = \max(x, 0.01x)$$

[63]

3.4.3

$$a(h_t) = \tanh(W_{hc}h_t + b_{hc})$$

Gulcehre

Temporal Automatic Relation Discovery In

104

c

Sequences TARDIS

[64]

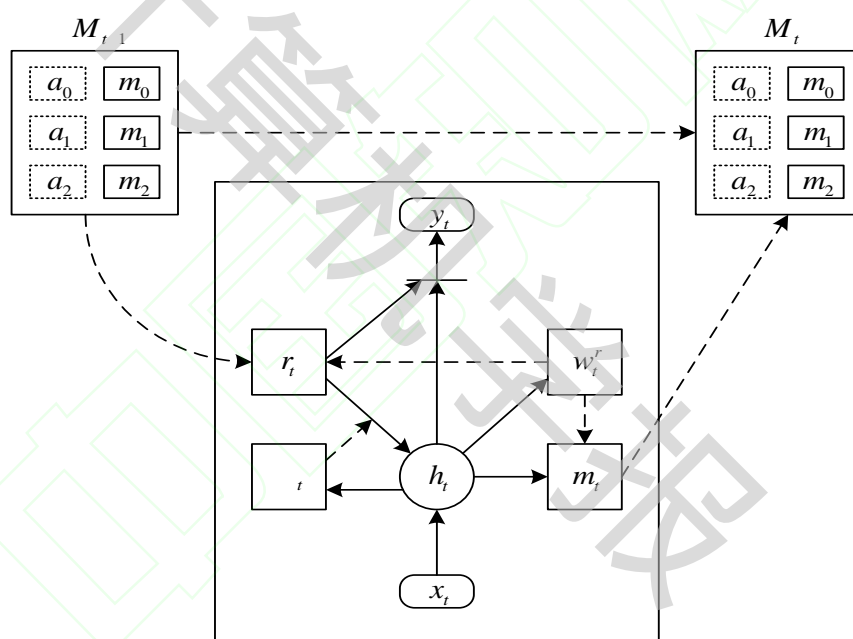
s

$$s = \text{LReLU}(W_{cs}c + b_{cs}), W_{cs} \in \mathbb{R}^{D \times D}, b_{cs} \in \mathbb{R}^D$$

121

TARDIS

y



29 TARDIS

TARDIS

29

h_t

M_t

RNN

r_t

x_t

w_t^r

M_t

$h_{t-1}, h_t, (x_t, h_{t-1}, r_t)$

$r_t, (M_t)^T w_t^r$ TARDIS

$M_t[i], W_m h_t$

TARDIS

$$\begin{aligned}
 & \text{TARDIS} & M_t & \text{LSTM} \\
 & A_t \hat{i} \quad i^{k' a} & C_t \hat{i} \quad i^{k' c} & c_t \\
 & M_t [A_t; C_t] \quad \square^{k(c a)} & & 3.4.4
 \end{aligned}$$

Sarthak Jain

[65]

$$\begin{aligned}
 & \bar{w}_t^r \\
 & i[i] \quad a^T \tanh(W_h h_t \quad W_x x_t \quad W_m M_t[i] \quad W_u u_t) \quad \text{Freebase} \\
 & \bar{w}_t^r \quad \text{soft max}(\quad) \quad n \quad n\text{-gram} \quad n \\
 & \{a, W_h^s, W_x^s, W_m^s, W_u^s\} \quad 123 \quad L \\
 & w_t^r \quad \bar{w}_t^r \quad \arg \max \bar{w}_t^r \\
 & u_t \quad 1 \quad F = \{f_1, L, f_{|F|}\} \\
 & u_t = \text{norm}(\overset{t-1}{\underset{i=1}{\mathbf{a}}} w_i^r) \quad 124 \quad f \quad (s, R, o) \quad R \quad s \quad o \\
 & u_t \quad 2 \quad q \hat{i} \quad i^d \\
 & \text{TARDIS} \quad \text{TARDIS} \quad \text{LSTM} \quad f \hat{i} \quad F \\
 & \text{Gumbel-sigmoid} \quad g(f) \quad h(f) \\
 & \text{LSTM} \quad c_t \quad g(f) \\
 & \tilde{c}_t \quad \tanh(\quad W_h^s h_t \quad W_x^s x_t \quad W_r^s r_t) \quad q \quad (s, R) \quad g(f) \quad q^T(s \quad R) \\
 & c_t \quad f_t c_{t-1} \quad i_t \tilde{c}_t \quad 125 \\
 & h_t \quad o_t \tanh(c_t) \quad \text{softmax} \quad h(f) \\
 & f_t, i_t, o_t \quad \text{LSTM} \quad h(f) \\
 & a_t, b_t
 \end{aligned}$$

$$F = \{f_1, L, f_{|F|}\}$$

q

Memory SAM ^[66]

N

SAM

1

$(\log N)$

(N)

(1)

2

$$f = (s, R, o)$$

SAM

1

q

$q \quad q$

$$h(f)(s \quad R)$$

126

3

f

$$\tilde{w}_t^R \quad \tilde{w}_t^R(s_i)M_t(s_i)$$

128

o

f

F

\tilde{w}_t^R

s_1, L, s_K

q

w_t^R

K

$0 \quad K$

$q \quad K$

$h(f)$

Approximate nearest neighbors ANN

\tilde{w}_t^R

2

$$L_{QA} = \sum_{(q,A) n=1}^L \frac{n}{L} \left\| \left\| F_n \right\|_{a \quad A} - \left\| A \right\|_{f \quad F_n} h(f) \right\|_o^2$$

127

D

$$w_t^W \quad (w_t^R \quad (1 \quad)I_t^U)$$

129

q

A

n

n

L

w_{t-1}^R

o

I_t^U

0

a

$h(f)$

K

I_t^U

3.4.5

U

1

0

$$U_T^{(1)}(i) = \sum_{t=0}^T l^{T-t} (w_t^W(i) + w_t^R(i)) \quad 130$$

r_t LSTM

l y_t

$$U_T^{(2)}(i) = T - \max \{t : w_t^W(i) + w_t^R(i) > d\} \quad 3.4.6$$

NTM

131

d LSTM Zhang

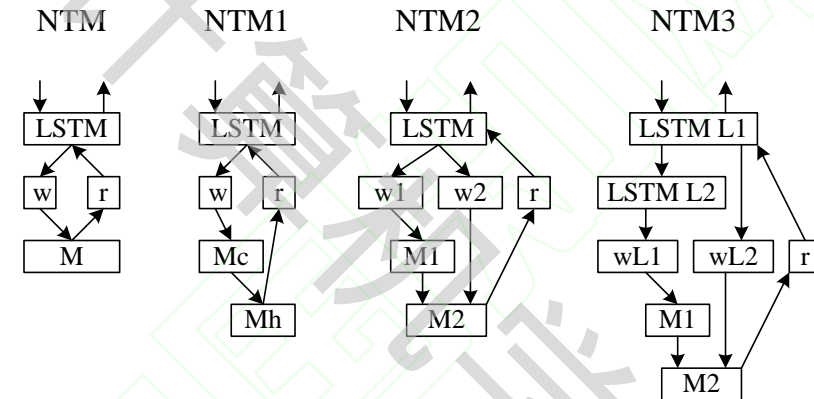
3

LSTM x_t

[67]

r_{t-1}

$$p_t = (q_t, a_t, a_t, g_t)$$



NTM

NTM1

Zhang

NTM1

NTM1

30

NTM3

30 NTM

1 NTM1

NTM

M_h M_h

M_h

t

M_c

M_h

$r(t)$

$c(t)$

h

M_c

a b

$M_h(t)$

t NTM1

2 NTM2

M_1 M_2

$$M_c(t) = h(M_c(t-1), w(t-1), c(t))$$

$$M_h(t) = aM_h(t-1) + bM_c(t) \quad 132$$

M_2

w_2

$$r(t) = w_r(t)M_h(t)$$

M_c $t-1$

$w(t-1)$

M_1

w_1

t

M_2 M_1 w_2 NTM2

$$M_1(t), w_1(t) = h(M_1(t-1), w_1(t-1), c(t))$$

$$\dot{M}_2(t), w_2(t) = h(M_2(t-1), w_2(t-1), c(t)) \quad 133$$

$$M_2(t) = a\dot{M}_2(t) + bM_1(t)$$

$$r(t) = w_r(t)M_2(t)$$

M_1 M_2

M_2

3 NTM3

[71] Luyang Li

[39] Jianpeng Cheng LSTM

[78] Tang D LSTM

[79] Chen P

Machine Comprehension

Pan LSTM [73]

Tran N LSTM

Full-orientation matching

4.2

Hamid Palangi LSTM [74]

Xu Jia LSTM [82]

Web

Arnav Kumar Jain [54]

Sun

J [80] Chunseong Park C Context Sequence Memory Network CSMN [83]

Mingxuan Wang RNN [75]

Instagram

Ganhotra [76]

Oriol Vinyals [40]

Adam Santoro [32]

Huang RNN - [77]

Wang J [88]

Donahue J
[87] Long-Term Recurrent Convolutional Networks
LTRCN LSTM

4.3

Kim Y B

[93] Speaker-Sensitive Dual Memory
Networks SSDMN
multi-turn slot tagging

Caiming Xiong

[43] Ma C [84]

Kim K M

AI

Huang C Z A

[94]

Joel Moniz

[86]

LSTM

text-to-speech TTS Tacotron2 TTS Naihan Li [95]

Scene labeling

Abdulnabi A H

[89]

CNN

4.4

Parisotto E [96]

Kaiser Ł
life-long

[90]

one-shot

2D

Baskar M K

[97]

Bornschein J

Parmar N

[98]

[91]

ImageNet

[99]

Fernando T

Conditional Generative Adversarial Networks, (Memory Augmented Networks, MC-GAN) [92]

Trang Pham

DeepCare^[100]

Pham T

[108]

Prakash A

[101] Condensed Memory

Neural Networks CMNN

Knowledge Tracing

Schwaller P

Zhang J

[109] Dynamic Key-Value Memory Networks

[102] Molecular DKV-MemNN

Transformer MT

Jonathan Woodbridge

LSTM

Domain Generate Algorithm DGA

[103]

DGA

Sprechmann P

DGA

Lee S W

[104]

4.5

Fernando T

[105]

9

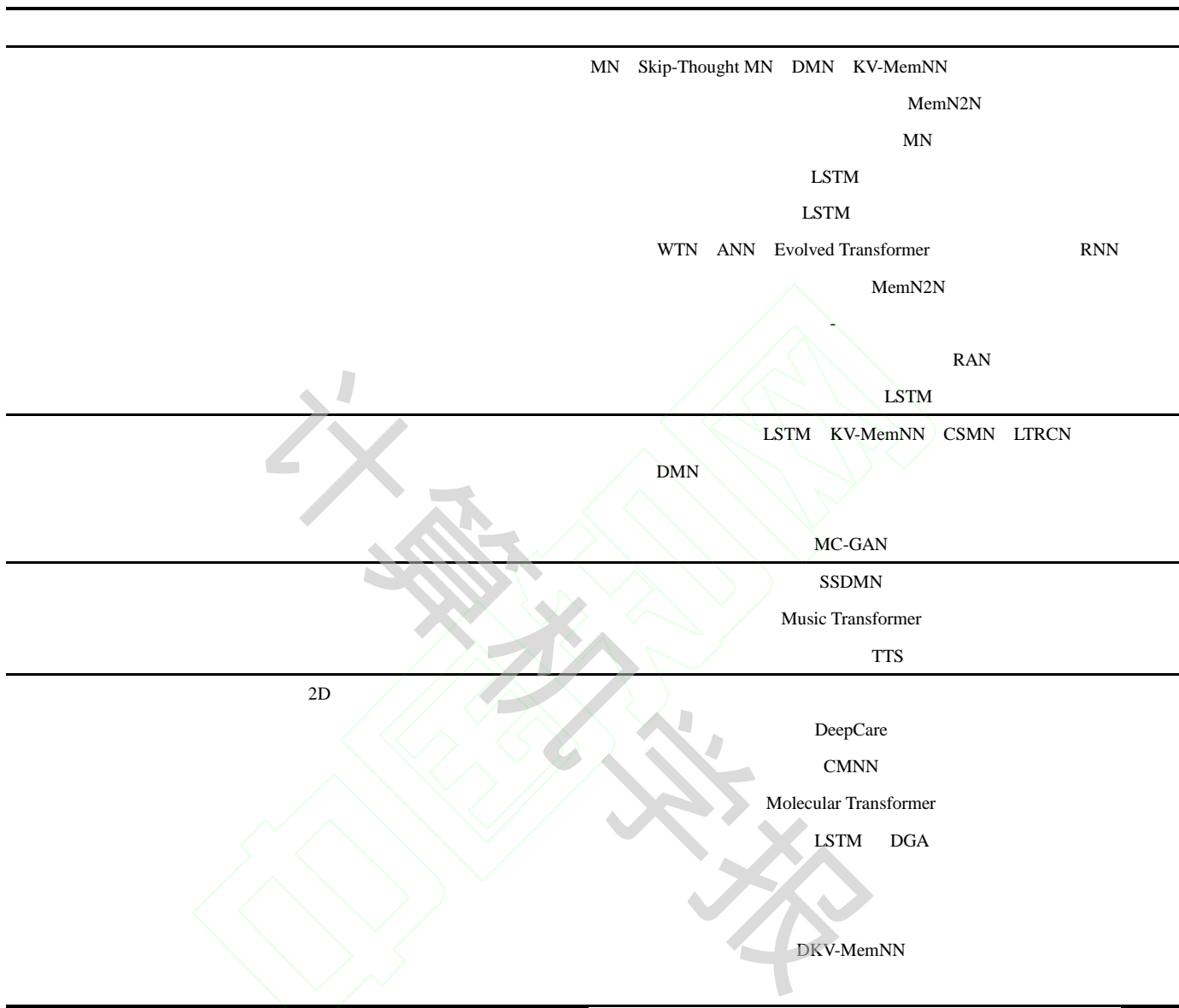
Zadeh A

[106]

2D

Khan A

[107]



5

5.1

RNN

LSTM

NTM

MN

RNN

MN
[111]

4

- - Tishby [112]

[113]

5

5.2

MN

LSTM

6

[23]

1

7

2

3 MN

MN

LSTM

MN

9

16

10

CNN

GAN [90]

11

12

13

[100-102]

[103]

[107]

[1] Sepp Hochreiter, Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation* 9(8): 1735-1780 (1997)

[2] Jason Weston, Sumit Chopra, Antoine Bordes. Memory Networks. *CoRR* abs/1410.3916 (2014).

[3] Alex Graves, Greg Wayne, Ivo Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.

[4] Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, Rob Fergus. End-to-end memory networks// Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, Roman Garnett. *Advances in neural information processing systems*. Montreal, Quebec, Canada. NIPS 2015: 2440-2448.

[5] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need[C]//*Advances in neural information processing systems*. 2017: 5998-6008.

[6] Klaus Greff, Rupesh Kumar Srivastava, Jan Koutn k, Bas R. Steunebrink, Jürgen Schmidhuber. LSTM: A search space odyssey// *IEEE Trans. Neural Netw. Learning Syst.* 28(10): 2222-2232 (2017)

[7] Fader A, Zettlemoyer L, Etzioni O. Paraphrase-driven learning for open question answering[C]//*Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2013: 1608-1618.

[8] Bordes A, Weston J, Usunier N. Open question answering with weakly supervised embedding models[C]//*Joint European conference*

14

15

on machine learning and knowledge discovery in databases. Springer, Berlin, Heidelberg, 2014: 165-180.

- [9] Mikolov T, Karafiát M, Burget L, et al. Recurrent neural network based language model[C]//Eleventh annual conference of the international speech communication association. 2010.
- [10] Kingma D P, Ba J. Adam: A method for stochastic optimization[J]. arXiv preprint arXiv:1412.6980, 2014.
- [11] Srivastava N, Hinton G, Krizhevsky A, et al. Dropout: a simple way to prevent neural networks from overfitting[J]. The journal of machine learning research, 2014, 15(1): 1929-1958.
- [12] Szegedy C, Vanhoucke V, Ioffe S, et al. Rethinking the inception architecture for computer vision[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2016: 2818-2826.
- [13] Kalchbrenner N, Espeholt L, Simonyan K, et al. Neural machine translation in linear time[J]. arXiv preprint arXiv:1610.10099, 2016.
- [14] Zhou J, Cao Y, Wang X, et al. Deep recurrent models with fast-forward connections for neural machine translation[J]. Transactions of the Association for Computational Linguistics, 2016, 4: 371-383.
- [15] Wu Y, Schuster M, Chen Z, et al. Google's neural machine translation system: Bridging the gap between human and machine translation[J]. arXiv preprint arXiv:1609.08144, 2016.
- [16] Gehring J, Auli M, Grangier D, et al. Convolutional sequence to sequence learning[C]//Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017: 1243-1252.
- [17] Shazeer N, Mirhoseini A, Maziarz K, et al. Outrageously large neural networks: The sparsely-gated mixture-of

- Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, Roman Garnett, Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, Barcelona, Spain, NIPS 2016: 3630-3638
- [41] Fei Liu, Julien Perez. Gated end-to-end memory networks// Mirella Lapata, Phil Blunsom, Alexander Koller, Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, Valencia, Spain, EACL 2017: 1-10
- [42] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, Richard Socher. Ask Me Anything: Dynamic Memory Networks for Natural Language Processing// Maria-Florina Balcan, Kilian Q. Weinberger, Proceedings of the 33rd International Conference on Machine Learning, New York City, NY, USA, ICML 2016: 1378-1387
- [43] Alexander H. Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, Jason Weston. Key-Value Memory Networks for Directly Reading Documents// Jian Su, Xavier Carreras, Kevin Duh, Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Austin, Texas, USA, EMNLP 2016: 1400-1409
- [44] Sarath Chandar, Sungjin Ahn, Hugo Larochelle, Pascal Vincent, Gerald Tesauro, Yoshua Bengio. Hierarchical memory networks. arXiv preprint arXiv:1605.07427, 2016.
- [45] Alex Auvolat, Pascal Vincent. Clustering is efficient for approximate maximum inner product search. arXiv preprint arXiv:1507.05910, 2015.
- [46] Zhong S. Efficient online spherical k-means clustering// Neural Networks, 2005. IJCNN'05. Proceedings. 2005 IEEE International Joint Conference on. IEEE, 2005, 5: 3180-3185.
- [47] Marcin Andrychowicz, Karol Kurach. Learning efficient algorithms with hierarchical attentive memory. arXiv preprint arXiv:1602.03218, 2016.
- [48] Ronald J. Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. Machine Learning 8: 229-256 (1992)
- [49] Da-Rong Liu, Shun-Po Chuang, Hung-yi Lee. Attention-based Memory Selection Recurrent Network for Language Modeling. arXiv preprint arXiv:1611.08656, 2016.
- [50] Antoine Bordes, Nicolas Usunier, Jason Weston, 2015. arXiv preprint arXiv:1506.02076, 2015.

- [65] Sarthak Jain. Question Answering over Knowledge Base using Factual Memory Networks// Proceedings of the Student Research Workshop, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language

- [86] Joel Moniz, Christopher J. Pal. Convolutional residual memory networks[J]. arXiv preprint arXiv:1606.05262, 2016.
- [87] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, Trevor Darrell. Long-Term Recurrent Convolutional Networks for Visual Recognition and Description// IEEE Transactions on Pattern Analysis and Machine Intelligence. 39(4): 677-691 (2017)
- [88] Junbo Wang, Wei Wang, Yan Huang, Liang Wang, Tieniu Tan. M3: Multimodal Memory Modelling for Video Captioning// 2018 IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, CVPR 2018: 7512-7520
- [89] Abrar H. Abdulnabi, Bing Shuai, Stefan Winkler, Gang Wang. Episodic CAMN: Contextual Attention-Based Memory Networks with Iterative Feedback for Scene Labeling// 2017 IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, CVPR 2017: 6278-6287
- [90] Lukasz Kaiser, Ofir Nachum, Aurko Roy, Samy Bengio. Learning to remember rare events. arXiv preprint arXiv:1703.03129, 2017.
- [91] Parmar N, Vaswani A, Uszkoreit J, et al. Image transformer[J]. arXiv preprint arXiv:1802.05751, 2018.
- [92] Tharindu Fernando, Simon Denman, Sridha Sridharan, Clinton Fookes. Task Specific Visual Saliency Prediction with Memory Augmented Conditional Generative Adversarial Networks// 2018 IEEE Winter Conference on Applications of Computer Vision, Lake Tahoe, NV, USA, WACV 2018: 1539-1548
- [93] Young-Bum Kim, Sungjin Lee, Ruhi Sarikaya. Speaker-sensitive dual memory networks for multi-turn slot tagging// 2017 IEEE Automatic Speech Recognition and Understanding Workshop, Okinawa, Japan, ASRU 2017: 541-546
- [94] Huang C Z A, Vaswani A, Uszkoreit J, et al. Music transformer: Generating music with long-term structure[J]. 2018.
- [95] Li N, Liu S, Liu Y, et al. Close to human quality TTS with transformer[J]. arXiv preprint arXiv:1809.08895, 2018.
- [96] Emilio Parisotto, Ruslan Salakhutdinov. Neural map: Structured memory for deep reinforcement learning. arXiv preprint arXiv:1702.08360, 2017.
- [97] Murali Karthick Baskar, Martin Karafiát, Lukáš Burget, Karel Veselý, Frantisek Grézl, Jan Cernocký. Residual memory networks: Feed-forward approach to learn long-term temporal dependencies// 2017 IEEE International Conference on Acoustics, Speech and Signal Processing, New Orleans, LA, USA, ICASSP 2017: 4810-4814
- [98] Jörg Bornschein, Andriy Mnih, Daniel Zoran, Danilo Jimenez Rezende. Variational Memory Addressing in Generative Models// Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, Roman Garnett, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, Long Beach, CA, USA, NIPS 2017: 3923-3932
- [99] Diederik P. Kingma, Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- [100] Trang Pham, Truyen Tran, Dinh Q. Phung, Svetha Venkatesh. DeepCare: A Deep Dynamic Memory Model for Predictive Medicine// James Bailey, Latifur Khan, Takashi Washio, Gillian Dobbie, Joshua Zhexue Huang, Ruili Wang, Advances in Knowledge Discovery and Data Mining - 20th Pacific-Asia Conference, Auckland, New Zealand, PAKDD (2) 2016: 30-41
- [101] Aaditya Prakash, Siyuan Zhao, Sadid A. Hasan, Vivek V. Datla, Kathy Lee, Ashequl Qadir, Joey Liu, Oladimeji Farri. Condensed Memory Networks for Clinical Diagnostic Inferencing// Satinder P. Singh, Shaul Markovitch, Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, San Francisco, California, USA, AAAI 2017: 3274-3280
- [102] Schwaller P, Laino T, Gaudin T, et al. Molecular transformer for chemical reaction prediction and uncertainty estimation[J]. arXiv preprint arXiv:1811.02633, 2018
- [103] Jonathan Woodbridge, Hyrum S. Anderson, Anjum Ahuja, Daniel Grant. Predicting domain generation algorithms with long short-term memory networks. arXiv preprint arXiv:1611.00791, 2016.
- [104] Sang-Woo Lee, Chung-yeon Lee, Dong-Hyun Kwak, Jung-Woo Ha, Jeonghee Kim, Byoung-Tak Zhang. Dual-memory neural networks for modeling cognitive activities of humans via wearable sensors. Neural Networks 92: 17-28 (2017)
- [105] Tharindu Fernando, Simon Denman, Aaron McFadyen, Sridha Sridharan, Clinton Fookes. Tree Memory Networks for modelling long-term temporal dependencies. Neurocomputing 304: 64-81 (2018)
- [106] Zadeh A, Liang P P, Mazumder N, et al Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, Louis-Philippe Morency. Memory Fusion Network for Multi-view Sequential Learning// Sheila A. McIlraith, Kilian Q. Weinberger, Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, AAAI 2018: 5634-5641
- [107] Muhammad Asjad Khan, Hung Le, Kien Do, Truyen Tran, Aditya Ghose, Hoa Dam, Renuka Sindhgatta. Memory-Augmented Neural Networks for Predictive Process Analytics. CoRR abs/1802.00938 (2018)
- [108] Trang Pham, Truyen Tran, Svetha Venkatesh. Relational dynamic memory networks. CoRR abs/1808.04247 (2018)

- [109] Jiani Zhang, Xingjian Shi, Irwin King, Dit-Yan Yeung. Dynamic Key-Value Memory Networks for Knowledge Tracing// Rick Barrett, Rick Cummings, Eugene Agichtein, Evgeniy Gabrilovich, Proceedings of the 26th International Conference on World Wide Web, Perth, Australia, WWW 2017: 765-774
- [110] Pablo Sprechmann, Siddhant M. Jayakumar, Jack W. Rae, Alexander Pritzel, Adrià Puigdomènech Badia, Benigno Uria, Oriol Vinyals, Demis Hassabis, Razvan Pascanu, Charles Blundell. Memory-based parameter adaptation. arXiv preprint arXiv:1802.10542, 2018.
- [111] , , , . [J]. Journal of Software, 2017, 11: 2905-2924.
- [112] Tishby N, Zaslavsky N. Deep learning and the information bottleneck principle[C]//2015 IEEE Information Theory Workshop (ITW). IEEE, 2015: 1-5.
- [113] , . : [J]. , 2018, 1(11): 86.

Background

Deep memory network is a general term for neural network models with memory function, which is mainly to solve the prediction problem of sequence-dependent dependence, and can be predicted by memorizing the effective information learned before. Memory network usually have independent memory modules or other structures capable of memory function. The former stores important information in an independently readable and writable memory and reads it when needed; while the latter method usually modify the internal structure of the cell to retain the information that needs to be remembered.

Deep memory network have achieved unprecedented performance in a wide variety of different application areas. For example, image classification, face recognition, human-level concept learning, playing Atari games and AlphaGo.

Deep memory network combines the benefits of memory network and deep learning. On one hand, memory network has a wider scope of applicability since it can enhance the memory of the model. On the other hand, deep learning can extract a good representation at different levels of abstraction, which disentangles better the factors of variations underlying the data.

In this paper we aim to survey and place in dg AJT befd[(m)17(e)-10()JT