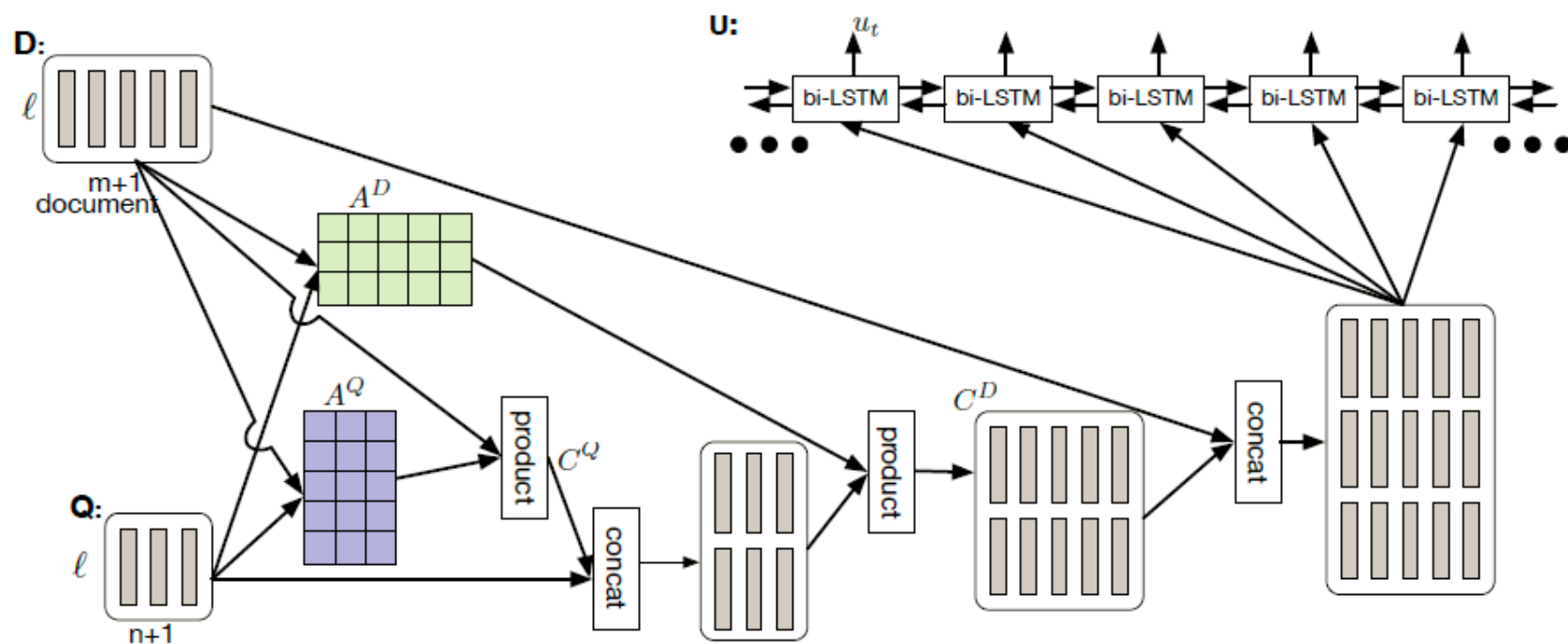


# DCN+

Caiming Xiong, Victor Zhong, Richard Socher

# Baseline DCN



## Baseline DCN

$$L = D^T Q \in R^{(m+1)*(n+1)}$$

$$A^Q = \textit{softmax}(L) \in R^{(m+1)*(n+1)} \text{ and}$$

$$A^D = \textit{softmax}(L^T) \in R^{(n+1)*(m+1)}$$

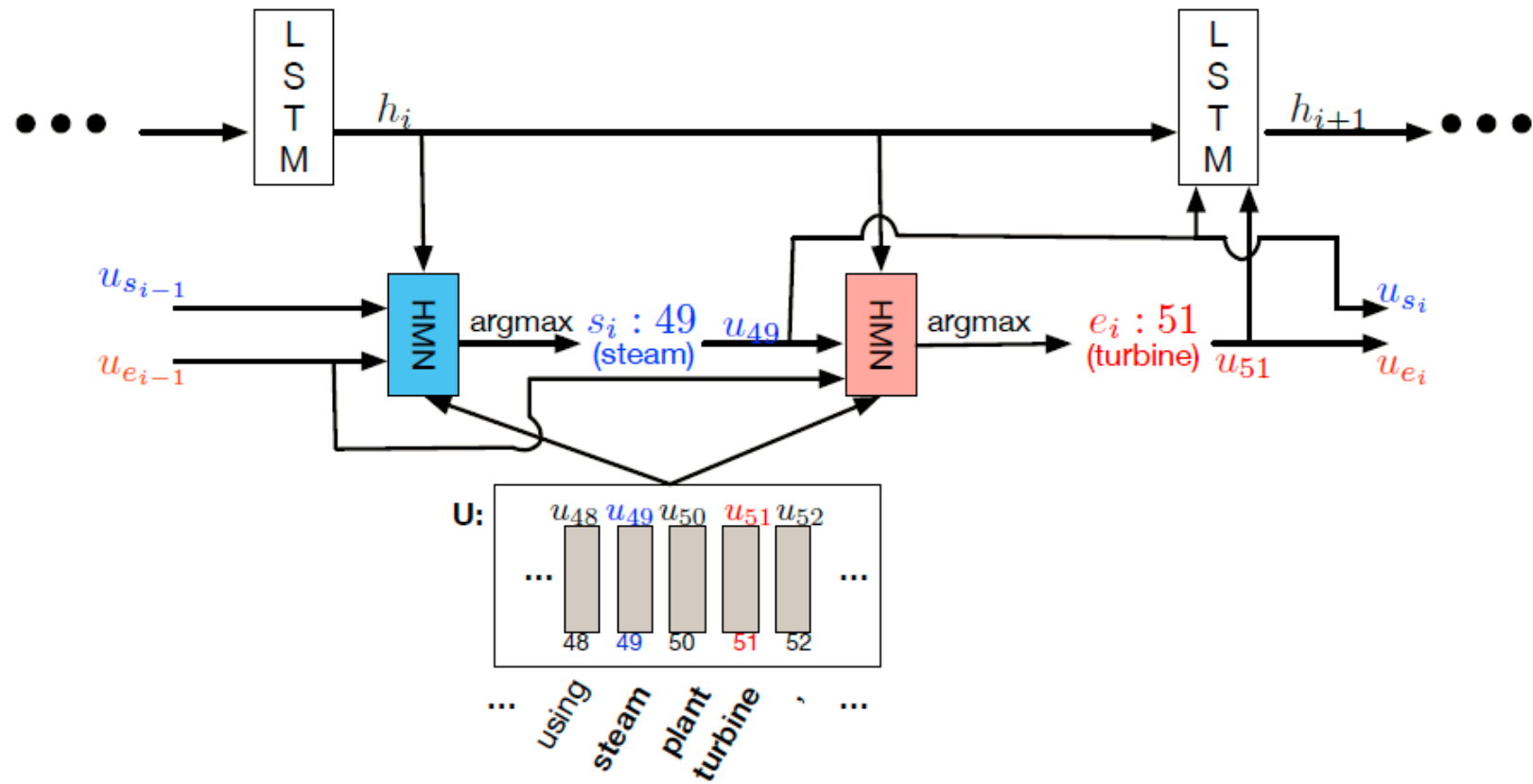
$$C^Q = DA^Q \in R^{l*(n+1)}$$

$$C^D = [Q; C^Q]A^D \in R^{2l*(m+1)}$$

$$u_t = \textit{BiLSTM}(u_{t-1}, u_{t+1}, [d_t; c_t^D]) \in R^{2l}$$

# Baseline DCN

## Dynamic Decoder



## Dynamic Decoder

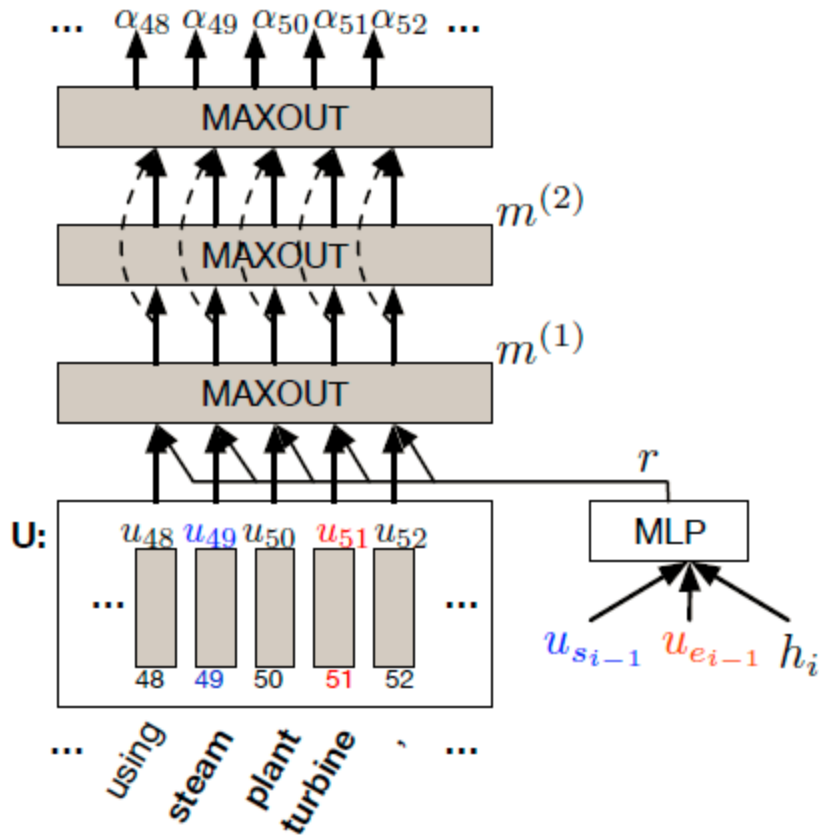
$$h_i = LSTM_{dec}(h_{i-1}, [u_{s_{i-1}}; u_{e_{i-1}}])$$

$$s_i = \operatorname{argmax}_t(\alpha_1, \dots, \alpha_m)$$

$$e_i = \operatorname{argmax}_t(\beta_1, \dots, \beta_m)$$

$$\alpha_t = HMN_{start}(u_t, h_i, u_{s_{i-1}}, u_{e_{i-1}})$$

# Dynamic Deocder



$$\begin{aligned}
 \text{HMN}(u_t, h_i, u_{s_{i-1}}, u_{e_{i-1}}) &= \max \left( W^{(3)} \begin{bmatrix} m_t^{(1)}; m_t^{(2)} \end{bmatrix} + b^{(3)} \right) \\
 r &= \tanh \left( W^{(D)} \begin{bmatrix} h_i; u_{s_{i-1}}; u_{e_{i-1}} \end{bmatrix} \right) \\
 m_t^{(1)} &= \max \left( W^{(1)} [u_t; r] + b^{(1)} \right) \\
 m_t^{(2)} &= \max \left( W^{(2)} m_t^{(1)} + b^{(2)} \right)
 \end{aligned}$$

# DCN+

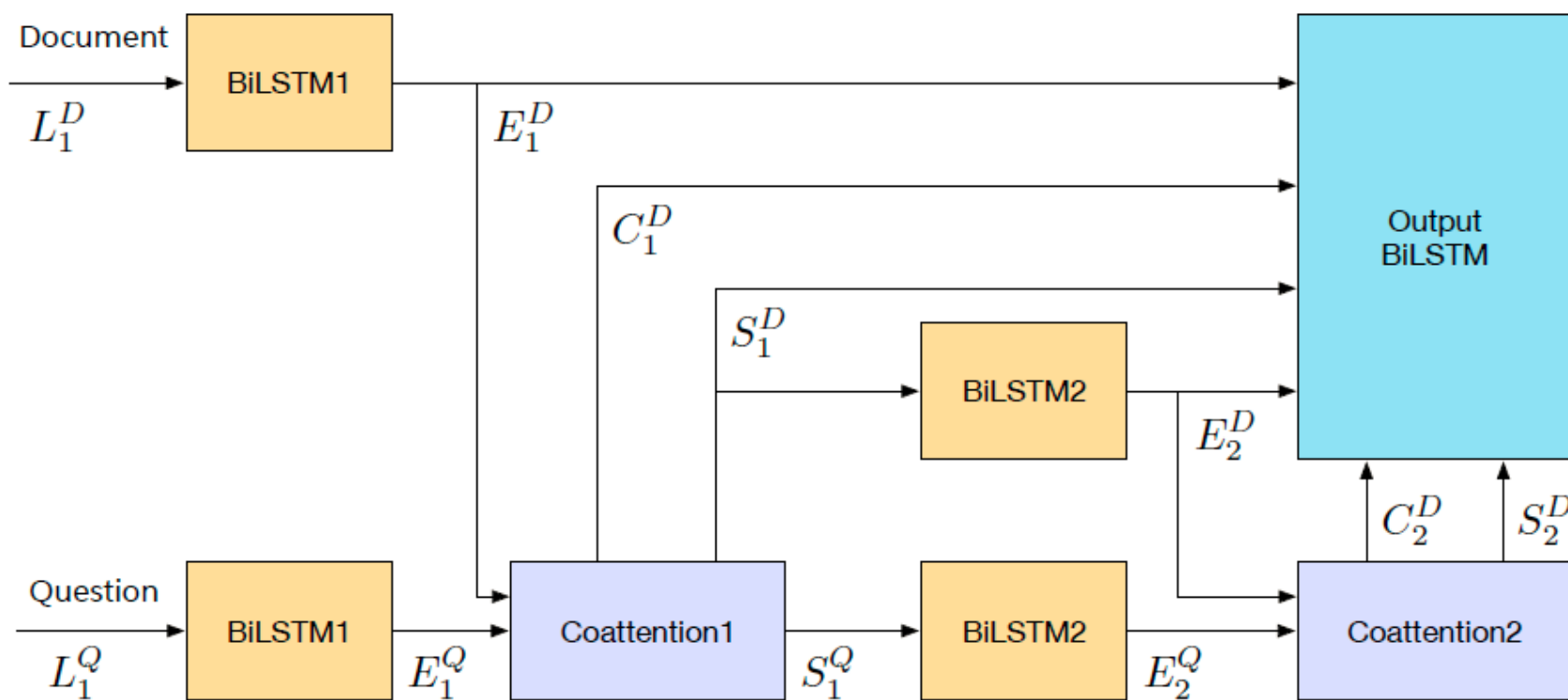


Figure 1: Deep residual coattention encoder.

First, encoding for D and Q:

$$E_1^D = biLSTM_1(L^D) \in R^{h*(m+1)}$$

$$E_1^Q = tanh(W biLSTM_1(L^Q) + b) \in R^{h*(n+1)}$$

Secondly, counting the representation  $Q2D_1$  and  $D2Q_1$ :

$$A = (E_1^Q)^T E_1^D \in R^{(m+1)*(n+1)}$$

$$S_1^D = E_1^D softmax(A^T) \in R^{h*(m+1)}$$

$$S_1^Q = E_1^Q softmax(A) \in R^{h*(n+1)}$$

Thirdly, counting the  $(D2Q_1)2D$  representation:

$$C_1^D = S_1^Q softmax(A^T) \in R^{h*m}$$

Fourthly, counting the second LSTM encoding for  $S_1^D$  and  $S_1^Q$ :

$$E_2^D = biLSTM_2(S_1^D) \in R^{2h*m}$$

$$E_2^Q = biLSTM_2(S_1^Q) \in R^{2h*n}$$

Last, doing the same coattention operation for  $E_2^D$  and  $E_2^Q$ .



## DCN+

In summary,

$$coattn_1(E_1^D, E_1^Q) \rightarrow S_1^D, S_1^Q, C_1^D$$

$$coattn_2(E_2^D, E_2^Q) \rightarrow S_2^D, S_2^Q, C_2^D$$

The output of encoder is:

$$U = biLSTM(concat(E_1^D; E_2^D; S_1^D; S_2^D; C_1^D; C_2^D)) \in R^{2h*m}$$

# Optimization Objective

$$l_{ce}(\Theta) = - \sum_t (\log p_t^{\text{start}}(s \mid s_{t-1}, e_{t-1}; \Theta) + \log p_t^{\text{end}}(e \mid s_{t-1}, e_{t-1}; \Theta))$$

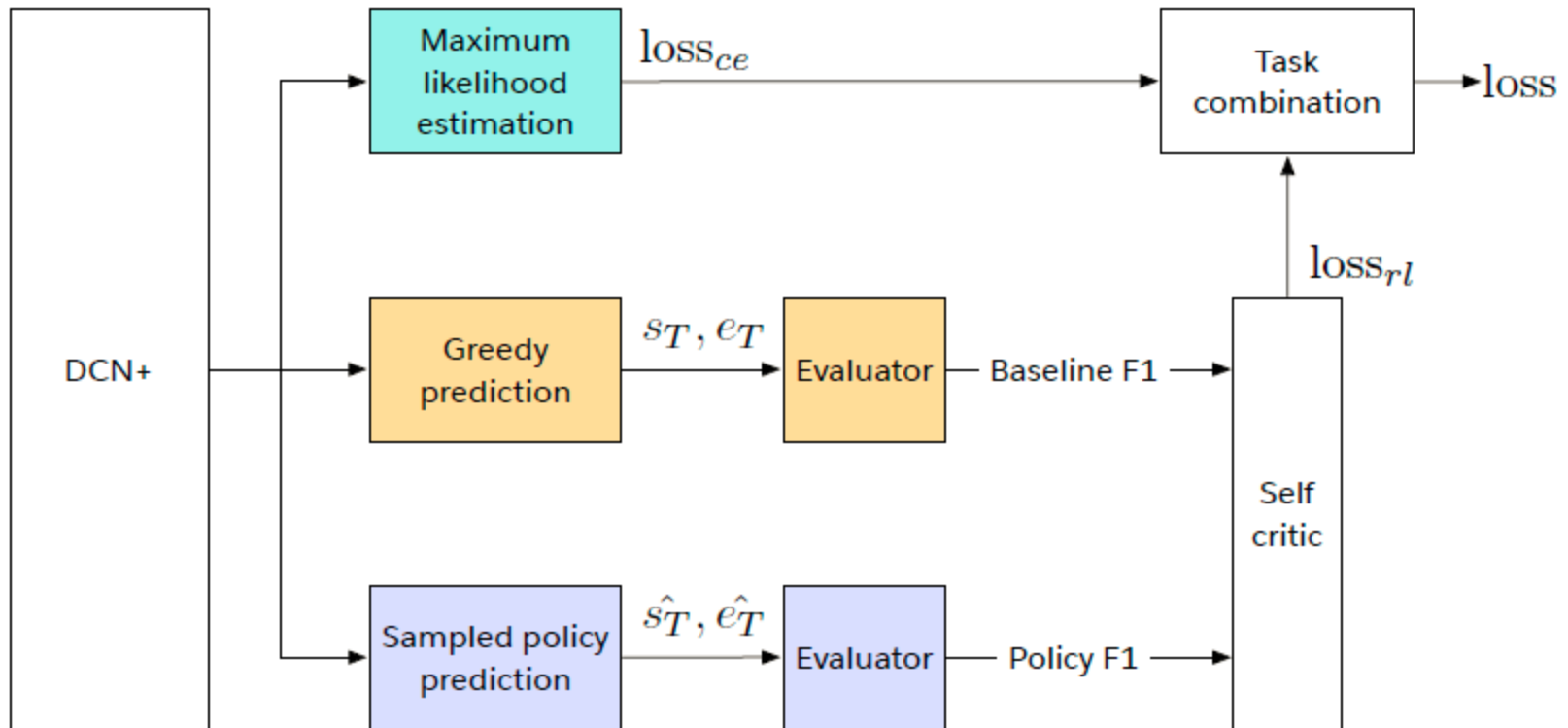
$$\begin{aligned} l_{rl}(\Theta) &= -\mathbb{E}_{\hat{\tau} \sim p_{\tau}} [R(s, e, \hat{s}_T, \hat{e}_T; \Theta)] \\ &\approx -\mathbb{E}_{\hat{\tau} \sim p_{\tau}} [F_1(\text{ans}(\hat{s}_T, \hat{e}_T), \text{ans}(s, e)) - F_1(\text{ans}(s_T, e_T), \text{ans}(s, e))] \end{aligned}$$

$$\nabla_{\Theta} l_{rl}(\Theta) = -\nabla_{\Theta} (\mathbb{E}_{\hat{\tau} \sim p_{\tau}} [R]) \quad (14)$$

$$= -\mathbb{E}_{\hat{\tau} \sim p_{\tau}} [R \nabla_{\Theta} \log p_{\tau}(\tau; \Theta)] \quad (15)$$

$$\begin{aligned} &= -\mathbb{E}_{\hat{\tau} \sim p_{\tau}} \left[ R \nabla_{\Theta} \left( \sum_t^T (\log p_t^{\text{start}}(\hat{s}_t | \hat{s}_{t-1}, \hat{e}_{t-1}; \Theta) + \log p_t^{\text{end}}(\hat{e}_t | \hat{s}_{t-1}, \hat{e}_{t-1}; \Theta)) \right) \right] \\ &\approx -R \nabla_{\Theta} \left( \sum_t^T (\log p_t^{\text{start}}(\hat{s}_t | \hat{s}_{t-1}, \hat{e}_{t-1}; \Theta) + \log p_t^{\text{end}}(\hat{e}_t | \hat{s}_{t-1}, \hat{e}_{t-1}; \Theta)) \right) \quad (16) \end{aligned}$$

# Optimization Objective



# Experiment

Model	EM	$\Delta$ EM	F1	$\Delta$ F1
DCN+ (ours)	74.5%	–	83.1%	–
- Deep residual coattention	73.1%	-1.4%	81.5%	-1.6%
- Mixed objective	73.8%	-0.7%	82.1%	-1.0%
- Mixture of experts	74.0%	-0.5%	82.4%	-0.7%
DCN w/ CoVe (baseline)	71.3%	-3.2%	79.9%	-3.2%

ps.CoVe and residual coattention is important.

Cross-entropy is important for RL.