

input word indices. layers, embedding.

↓ embedding layer

$$\text{embeddings: } C_1, \dots, C_N \in \mathbb{R}^D \quad \left. \begin{array}{l} q_1, \dots, q_m \in \mathbb{R}^D \end{array} \right\} V_i$$

↓ projection

$$h_i = W_{\text{proj}} \cdot V_i \in \mathbb{R}^H$$

$\mathbb{R}^{H \times D}$

↓ Highway Network x 2

gate →  $g = \sigma(W_g h_i + b_g) \in \mathbb{R}^H$

transform →  $t = \text{ReLU}(W_t h_i + b_t) \in \mathbb{R}^H$

$h_i' = g \odot t + (1-g) \odot h_i \in \mathbb{R}^H$

↓ Encoder layer. layers, RNN Encoder.

↓ bidirectional LSTM  
(1 layer)

$$h_{i,\text{fwd}} = \text{LSTM}(h_{i-1}, h_i) \in \mathbb{R}^H$$

$$h_{i,\text{rev}} = \text{LSTM}(h_{i+1}, h_i) \in \mathbb{R}^H$$

$$h_i' = [h_{i,\text{fwd}}; h_{i,\text{rev}}] \in \mathbb{R}^{2H}$$

↓ Attention Layer

$$C_1, \dots, C_N \in \mathbb{R}^{2H}$$

$$q_1, \dots, q_m \in \mathbb{R}^{2H}$$

$$S_{ij} = W_{\text{attn}} [C_i; q_j] \in \mathbb{R}$$

$\mathbb{R}^{6H}$  weight matrix

C2Q {

$$\bar{S}_{i,:} = \text{softmax}(S_{i,:}) \in \mathbb{R}^M \quad \forall i \in \{1, \dots, N\}$$

attention distr.

$$a_i = \sum_{j=1}^M \bar{S}_{i,j} q_j \in \mathbb{R}^{2H} \quad \forall i \in \{1, \dots, N\}$$

weighted sums of question hidden states.  
attention output

Q2C {

$$\bar{S}_{:,j} = \text{softmax}(S_{:,j}) \in \mathbb{R}^N \quad \forall j \in \{1, \dots, M\}$$

$$s' = \bar{S} \cdot \bar{S}^T \in \mathbb{R}^{N \times N}$$

$$b_i = \sum_{j=1}^M \bar{S}_{i,j} C_j \in \mathbb{R}^{2H} \quad \forall i \in \{1, \dots, N\}$$

output ←  $g_i = [C_i; a_i; C_i \cdot a_i; C_i \cdot b_i] \in \mathbb{R}^{8H} \quad \forall i \in \{1, \dots, N\}$

layers, RNN Encoder.

Modeling layer.

bidirectional LSTM (> layers)

$$m_{i,\text{fwd}} = \text{LSTM}(m_{i-1}, q_i) \in \mathbb{R}^H$$

$$m_{i,\text{rev}} = \text{LSTM}(m_{i+1}, q_i) \in \mathbb{R}^H$$

$$m_i = [m_{i,\text{fwd}}; m_{i,\text{rev}}] \in \mathbb{R}^{2H}$$

layers, BiDA Output.

Output layer

bidirectional LSTM

$$m_{i',\text{fwd}} = \text{LSTM}(m_{i'-1}, m_i) \in \mathbb{R}^H$$

$$m_{i',\text{rev}} = \text{LSTM}(m_{i'+1}, m_i) \in \mathbb{R}^H$$

$$m_{i'} = [m_{i',\text{fwd}}; m_{i',\text{rev}}] \in \mathbb{R}^{2H}$$

$$G \in \mathbb{R}^{8H \times N} \quad [g_1, \dots, g_N]$$

$$M, M' \in \mathbb{R}^{2H \times N} \quad \begin{bmatrix} m_1 & \dots & m_N \\ m'_1 & \dots & m'_N \end{bmatrix}$$

log-scale loss fn

$$P_{\text{start}} = \text{softmax}(W_{\text{start}} [G; M])$$

$$P_{\text{end}} = \text{softmax}(W_{\text{end}} [G; M'])$$

$\mathbb{R}^{1 \times 10H}$