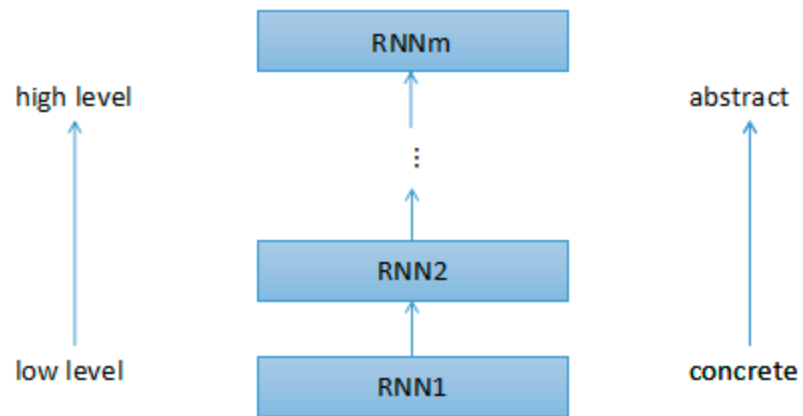# FusionNet

# Motivation

1.The encoding for word of the whole history procedure is important.
2.The fully-aware multi-level attention fusion is important.

# Contributions

This paper proposes a novel attention mechanism with following three contributions:

1.The concept of "history of words" to build the attention using complete information from lowest word-level embedding up to the highest semantic-level representation.

2.It proposes a novel attention scoring function.

3.It proposes a fully-aware multi-level fusion to exploit information layer by layer.

# Multi-Level

# Scoring Function

$$S_{ij} = f(U(\mathrm{HoW}_i^A))^T D \, f(U(\mathrm{HoW}_j^B))$$

# Fully-Aware Fusion Network

## Input Vectors

context: 921-dim {300-dim GloVe embedding, 600-dim contextualized vector, 12-dim POS embedding, 8-dim NER embedding, 1-dim normalized term frequency}
question: 900-dim {300-dim GloVe embedding, 600-dim conextualized vector}

$$C : \{w_1^C, ..., w_m^C\} \in R^{900+20+1}$$
$$Q : \{w_1^Q, ..., w_n^Q\} \in R^{900}$$

# Fully-aware Multi-level Fusion: Word-level

$$\hat{g}_i^C = \sum_j \alpha_{ij}\, g_j^Q, \quad \alpha_{ij} \propto \exp(S(g_i^C, g_j^Q)), \quad S(x,y) = \mathrm{ReLU}(Wx)^T \mathrm{ReLU}(Wy)$$

$em_i$ is created for each word in C to indicate wheter the word occurs in the question Q.

$$\tilde{w}_i^C = [w_i^C; \mathrm{em}_i; \hat{g}_i^C]$$

# Reading

$$h_1^{Cl}, \ldots, h_m^{Cl} = \mathrm{BiLSTM}(\tilde{w}_1^C, \ldots, \tilde{w}_m^C), \quad h_1^{Ql}, \ldots, h_n^{Ql} = \mathrm{BiLSTM}(w_1^Q, \ldots, w_n^Q),$$

$$h_1^{Ch}, \ldots, h_m^{Ch} = \mathrm{BiLSTM}(h_1^{Cl}, \ldots, h_m^{Cl}), \quad h_1^{Qh}, \ldots, h_n^{Qh} = \mathrm{BiLSTM}(h_1^{Ql}, \ldots, h_n^{Ql}).$$

Hence low-level and high-level concept $h^l, h^h \in \mathbb{R}^{250}$ are created for each word.

# Question Understanding

$$U_Q = \{u_1^Q, \ldots, u_n^Q\} = \text{BiLSTM}([h_1^{Ql}; h_1^{Qh}], \ldots, [h_n^{Ql}; h_n^{Qh}]).$$

where $\{u_i^Q \in \mathbb{R}^{250}\}_{i=1}^n$ are the understanding vectors for $Q$.

# Fully-aware Multi-level Fusion: Higher-level

$$\text{HoW}_i^C = [g_i^C; c_i^C; h_i^{Cl}; h_i^{Ch}], \quad \text{HoW}_i^Q = [g_i^Q; c_i^Q; h_i^{Ql}; h_i^{Qh}] \in \mathbb{R}^{1400},$$

1. Low-level fusion: $\hat{h}_i^{Cl} = \sum_j \alpha_{ij}^l h_j^{Ql}, \quad \alpha_{ij}^l \propto \exp(S^l(\text{HoW}_i^C, \text{HoW}_j^Q)).$

2. High-level fusion: $\hat{h}_i^{Ch} = \sum_j \alpha_{ij}^h h_j^{Qh}, \quad \alpha_{ij}^h \propto \exp(S^h(\text{HoW}_i^C, \text{HoW}_j^Q)).$

3. Understanding fusion: $\hat{u}_i^C = \sum_j \alpha_{ij}^u u_j^Q, \quad \alpha_{ij}^u \propto \exp(S^u(\text{HoW}_i^C, \text{HoW}_j^Q)).$

$$V_C = \{v_1^C, \ldots, v_m^C\} = \text{BiLSTM}([h_1^{Cl}; h_1^{Ch}; \hat{h}_1^{Cl}; \hat{h}_1^{Ch}; \hat{u}_1^C], \ldots, [h_m^{Cl}; h_m^{Ch}; \hat{h}_m^{Cl}; \hat{h}_m^{Ch}; \hat{u}_m^C]).$$

# Fully-aware Self-boosted Fusion

$$\text{HoW}_i^C = [\boldsymbol{g}_i^C; \boldsymbol{c}_i^C; \boldsymbol{h}_i^{Cl}; \boldsymbol{h}_i^{Ch}; \hat{\boldsymbol{h}}_i^{Cl}; \hat{\boldsymbol{h}}_i^{Ch}; \hat{\boldsymbol{u}}_i^C; \boldsymbol{v}_i^C] \in \mathbb{R}^{2400}.$$

Then perform fully-aware attention,

$$\hat{\boldsymbol{v}}_i^C = \sum_j \alpha_{ij}^s \boldsymbol{v}_j^C, \quad \alpha_{ij}^s \propto \exp(S^s(\text{HoW}_i^C, \text{HoW}_j^C)).$$

The final context representation is

$$U_C = \{\boldsymbol{u}_1^C, \ldots, \boldsymbol{u}_m^C\} = \text{BiLSTM}([\boldsymbol{v}_1^C; \hat{\boldsymbol{v}}_1^C], \ldots, [\boldsymbol{v}_m^C; \hat{\boldsymbol{v}}_m^C]).$$

# Output

Through the above operations, we can get

$$U_C = u_1^C, ..., u_{m'}^C$$
$$U_Q = u_1^Q, ..., u_n^Q$$

Then we get the vector representation of Q

$$u^Q = \sum_i \beta_i u_i^Q$$

For start,

$$P_i^S \propto \exp((\boldsymbol{u}^Q)^T W_S \boldsymbol{u}_i^C),$$

For end,

$$\boldsymbol{v}^Q = \text{GRU}(\boldsymbol{u}^Q, \sum_i P_i^{\check{S}} \boldsymbol{u}_i^C)$$

$$P_i^E \propto \exp((\boldsymbol{v}^Q)^T W_E \boldsymbol{u}_i^C)$$

# Experiments

| Attention Function | EM / F1 |
|---|---|
| Additive (MLP) | 71.8 / 80.1 |
| Multiplicative | 72.1 / 80.6 |
| Scaled Multiplicative | 72.4 / 80.7 |
| Scaled Multiplicative + ReLU | 72.6 / 80.8 |
| Symmetric Form | 73.1 / 81.5 |
| **Symmetric Form + ReLU** | **75.3 / 83.6** |

| Configuration | | Dev EM / F1 |
|---|---|---|
| $C, Q$ Fusion | Self $C$ | |
| High-level | None | 64.6 / 73.2 |
| FA High-level | None | 73.3 / 81.4 |
| FA All-level | None | 72.3 / 80.7 |
| FA Multi-level | None | 74.6 / 82.7 |
| FA Multi-level | Normal | 74.4 / 82.6 |
| FA Multi-level | FA | **75.3 / 83.6** |