

Appendix

A Implementation Details

A.1 MI-LSTM

Our MI-LSTM (without peephole connection) in experiments has the follow formulation:

$$\begin{aligned} \mathbf{z}_t &= \tanh(\boldsymbol{\alpha}_z \odot \mathbf{W}_z \mathbf{x}_t \odot \mathbf{U}_z \mathbf{h}_{t-1} + \boldsymbol{\beta}_{z,1} \odot \mathbf{U}_z \mathbf{h}_{t-1} + \boldsymbol{\beta}_{z,2} \odot \mathbf{W}_z \mathbf{x}_t + \mathbf{b}_z) \\ \mathbf{i}_t &= \sigma(\boldsymbol{\alpha}_i \odot \mathbf{W}_i \mathbf{x}_t \odot \mathbf{U}_i \mathbf{h}_{t-1} + \boldsymbol{\beta}_{i,1} \odot \mathbf{U}_i \mathbf{h}_{t-1} + \boldsymbol{\beta}_{i,2} \odot \mathbf{W}_i \mathbf{x}_t + \mathbf{b}_i) \\ \mathbf{f}_t &= \sigma(\boldsymbol{\alpha}_f \odot \mathbf{W}_f \mathbf{x}_t \odot \mathbf{U}_f \mathbf{h}_{t-1} + \boldsymbol{\beta}_{f,1} \odot \mathbf{U}_f \mathbf{h}_{t-1} + \boldsymbol{\beta}_{f,2} \odot \mathbf{W}_f \mathbf{x}_t + \mathbf{b}_f) \\ \mathbf{c}_t &= \mathbf{i}_t \odot \mathbf{z}_t + \mathbf{f}_t \odot \mathbf{c}_{t-1} \\ \mathbf{o}_t &= \sigma(\boldsymbol{\alpha}_o \odot \mathbf{W}_o \mathbf{x}_t \odot \mathbf{U}_o \mathbf{h}_{t-1} + \boldsymbol{\beta}_{o,1} \odot \mathbf{U}_o \mathbf{h}_{t-1} + \boldsymbol{\beta}_{o,2} \odot \mathbf{W}_o \mathbf{x}_t + \mathbf{b}_o) \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \end{aligned}$$

block input
input gate
forget gate
cell state
output gate
block output

where $\{\boldsymbol{\alpha}_*, \boldsymbol{\beta}_{*,1}, \boldsymbol{\beta}_{*,2}\}_{*=\{z,i,f,o\}}$ are bias vectors, σ denotes the sigmoid function.

A.2 MI-GRU

Our MI-GRU in experiments has the follow formulation:

$$\begin{aligned} \mathbf{z}_t &= \sigma(\boldsymbol{\alpha}_z \odot \mathbf{W}_z \mathbf{x}_t \odot \mathbf{U}_z \mathbf{h}_{t-1} + \boldsymbol{\beta}_{z,1} \odot \mathbf{U}_z \mathbf{h}_{t-1} + \boldsymbol{\beta}_{z,2} \odot \mathbf{W}_z \mathbf{x}_t + \mathbf{b}_z) && \text{update gate} \\ \mathbf{r}_t &= \sigma(\boldsymbol{\alpha}_r \odot \mathbf{W}_r \mathbf{x}_t \odot \mathbf{U}_r \mathbf{h}_{t-1} + \boldsymbol{\beta}_{r,1} \odot \mathbf{U}_r \mathbf{h}_{t-1} + \boldsymbol{\beta}_{r,2} \odot \mathbf{W}_r \mathbf{x}_t + \mathbf{b}_r) && \text{reset gate} \\ \tilde{\mathbf{h}}_t &= \tanh(\boldsymbol{\alpha}_h \odot \mathbf{W}_h \mathbf{x}_t \odot \mathbf{U}_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \boldsymbol{\beta}_{h,1} \odot \mathbf{U}_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \boldsymbol{\beta}_{h,2} \odot \mathbf{W}_h \mathbf{x}_t + \mathbf{b}_h) && \text{candidate activation} \\ \mathbf{h}_t &= (1 - \odot \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_{t-1} && \text{hidden activation} \end{aligned}$$

where $\{\boldsymbol{\alpha}_*, \boldsymbol{\beta}_{*,1}, \boldsymbol{\beta}_{*,2}\}_{*=\{z,r,h\}}$ are bias vectors, σ denotes the sigmoid function.