
Approximating Stacked and Bidirectional Recurrent Architectures with the Delayed Recurrent Neural Network: Supplementary Material

A. Theorem 1 Proof

Let us recall the notation introduced in the main paper. We use superscript (i) to refer to a weight matrix or vector related to layer i in a stacked network, e.g., $\mathbf{W}_h^{(i)}$, or $\mathbf{h}_t^{(i)}$. For a single-layer d-RNN, we refer to weight matrices and related vectors with "hat", e.g., $\hat{\mathbf{W}}_h$ or $\hat{\mathbf{h}}_t$. Additionally, we define the block notation as subvector $\hat{\mathbf{v}}_t^{\{i\}}$ refers to the i -th block of vector $\hat{\mathbf{v}}_t$ composed of k blocks. The blocks follow the definition in Equations (3)-(5).

Proof of Theorem 1. We prove Theorem 1 by induction on the sequence length t . First, we show that for $t = 1$ the stacked RNN and the d-RNN with the constrained weights are equivalent. Namely, for $t = 1$ we show that the outputs and the hidden states are the same, i.e. $\hat{\mathbf{y}}_k = \mathbf{y}_1$ and $\hat{\mathbf{h}}_i^{\{i\}} = \mathbf{h}_1^{(i)}$, respectively. Without loss of generality, we have for any i in $1 \dots k$ the following:

$$\begin{aligned}
 \hat{\mathbf{h}}_i^{\{i\}} &= f^{\{i\}} \left(\hat{\mathbf{W}}_x \mathbf{x}_i + \hat{\mathbf{W}}_h \hat{\mathbf{h}}_{i-1} + \hat{\mathbf{b}}_h \right) \\
 &= f \left(\hat{\mathbf{W}}_x^{\{i\}} \mathbf{x}_i + \mathbf{W}_x^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i-1\}} + \mathbf{W}_h^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i\}} + \mathbf{b}_h^{(i)} \right) \\
 &= f \left(\mathbf{0} + \mathbf{W}_x^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i-1\}} + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= f \left(\mathbf{W}_x^{(i)} \cdot \right. \\
 &\quad \left. f \left(\mathbf{W}_x^{(i-1)} \hat{\mathbf{h}}_{i-2}^{\{i-2\}} + \mathbf{W}_h^{(i-1)} \hat{\mathbf{h}}_{i-2}^{\{i-1\}} + \mathbf{b}_h^{(i-1)} \right) \right. \\
 &\quad \left. + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= f \left(\mathbf{W}_x^{(i)} \cdot \right. \\
 &\quad \left. f \left(\mathbf{W}_x^{(i-1)} \hat{\mathbf{h}}_{i-2}^{\{i-2\}} + \mathbf{W}_h^{(i-1)} \mathbf{h}_0^{(i-1)} + \mathbf{b}_h^{(i-1)} \right) \right. \\
 &\quad \left. + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= \dots \\
 &= f \left(\mathbf{W}_x^{(i)} \dots f \left(\mathbf{W}_x^{(j)} \dots \right. \right. \\
 &\quad \left. \left. f \left(\mathbf{W}_x^{(2)} \cdot f \left(\mathbf{W}_x^{(1)} \mathbf{x}_1 + \mathbf{W}_h^{(1)} \mathbf{h}_0^{(1)} + \mathbf{b}_h^{(1)} \right) \right. \right. \\
 &\quad \left. \left. + \mathbf{W}_h^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_h^{(2)} \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_h^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_h^{(j)} \right) \dots + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \left. \right)
 \end{aligned}$$

$$\begin{aligned}
 &= f \left(\mathbf{W}_x^{(i)} \dots f \left(\mathbf{W}_x^{(j)} \dots \right. \right. \\
 &\quad \left. \left. f \left(\mathbf{W}_x^{(2)} \mathbf{h}_1^{(1)} + \mathbf{W}_h^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_h^{(2)} \right) \right. \right. \\
 &\quad \left. \left. \dots + \mathbf{W}_h^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_h^{(j)} \right) \dots + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= \dots \\
 &= f \left(\mathbf{W}_x^{(i)} \dots \right. \\
 &\quad \left. f \left(\mathbf{W}_x^{(j)} \mathbf{h}_1^{(j-1)} + \mathbf{W}_h^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_h^{(j)} \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= \dots \\
 &= f \left(\mathbf{W}_x^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_h^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_h^{(i)} \right) \\
 &= \mathbf{h}_1^{(i)},
 \end{aligned}$$

where we used the initialization assumption $\hat{\mathbf{h}}_{i-1}^{\{i\}} = \mathbf{h}_0^{(i)}$ for all $i = 1 \dots k$, and the definition of the hidden state in Equations (3)-(4) for $j - 1$ blocks, in the previous steps. In particular, we have for $j = k$,

$$\hat{\mathbf{h}}_k^{\{k\}} = \mathbf{h}_1^{(k)}.$$

Plugging this result and the definition of the output weights and biases in Equation (8) into Equation (2) for computing the output, we obtain

$$\begin{aligned}
 \hat{\mathbf{y}}_k &= g \left(\hat{\mathbf{W}}_o \hat{\mathbf{h}}_k + \hat{\mathbf{b}}_o \right) \\
 &= g \left(\mathbf{W}_o \hat{\mathbf{h}}_k^{\{k\}} + \mathbf{b}_o \right) \\
 &= g \left(\mathbf{W}_o \mathbf{h}_1^{(k)} + \mathbf{b}_o \right) \\
 &= \mathbf{y}_1.
 \end{aligned} \tag{A.1}$$

Which concludes the basis of the induction.

Next, we assume that $\hat{\mathbf{h}}_{t+i-1}^{\{i\}} = \mathbf{h}_t^{(i)}$ for all $1 \leq i \leq k$ and $t \leq T - 1$, and prove that it holds for the hidden states for all layers when $t = T$: $\hat{\mathbf{h}}_{T+i-1}^{\{i\}} = \mathbf{h}_T^{(i)}$, $\forall 1 \leq i \leq k$. Without loss of generality, we have for the hidden state $\hat{\mathbf{h}}_{T+i-1}^{\{i\}}$ in

constrained weights single-layer d-RNN that,

$$\begin{aligned}
 \hat{\mathbf{h}}_{T+i-1}^{\{i\}} &= f^{\{i\}} \left(\hat{\mathbf{W}}_{\mathbf{x}} \mathbf{x}_{T+i-1} + \hat{\mathbf{W}}_{\mathbf{h}} \hat{\mathbf{h}}_{T+i-2} + \hat{\mathbf{b}}_{\mathbf{h}} \right) \\
 &= f \left(\hat{\mathbf{W}}_{\mathbf{x}}^{\{i\}} \mathbf{x}_{T+i-1} + \mathbf{W}_{\mathbf{x}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i-1\}} \right. \\
 &\quad \left. + \mathbf{W}_{\mathbf{h}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= f \left(\mathbf{0} + \mathbf{W}_{\mathbf{x}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i-1\}} + \mathbf{W}_{\mathbf{h}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \cdot f \left(\mathbf{W}_{\mathbf{x}}^{(i-1)} \hat{\mathbf{h}}_{T+i-3}^{\{i-2\}} + \right. \right. \\
 &\quad \left. \left. \mathbf{W}_{\mathbf{h}}^{(i-1)} \hat{\mathbf{h}}_{T+i-3}^{\{i-1\}} + \mathbf{b}_{\mathbf{h}}^{(i-1)} \right) \right. \\
 &\quad \left. + \mathbf{W}_{\mathbf{h}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= \dots \\
 &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \dots f \left(\mathbf{W}_{\mathbf{x}}^{(j)} \dots \right. \right. \\
 &\quad \left. \left. f \left(\mathbf{W}_{\mathbf{x}}^{(2)} f \left(\mathbf{W}_{\mathbf{x}}^{(1)} \mathbf{x}_T + \mathbf{W}_{\mathbf{h}}^{(1)} \hat{\mathbf{h}}_{T-1}^{\{1\}} + \mathbf{b}_{\mathbf{h}}^{(1)} \right) \right) \right. \right. \\
 &\quad \left. \left. + \mathbf{W}_{\mathbf{h}}^{(2)} \hat{\mathbf{h}}_{T-1}^{\{2\}} + \mathbf{b}_{\mathbf{h}}^{(2)} \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(j)} \hat{\mathbf{h}}_{T+j-2}^{\{j\}} + \mathbf{b}_{\mathbf{h}}^{(j)} \right) \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{h}}^{(i)} \right)
 \end{aligned}$$

From the inductive assumption we have $\hat{\mathbf{h}}_{T+j-2}^{\{j\}} = \mathbf{h}_{T-1}^{(j)}$ for all $1 \leq j \leq k$, then it follows

$$\begin{aligned}
 \hat{\mathbf{h}}_{T+i-1}^{\{i\}} &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \dots f \left(\mathbf{W}_{\mathbf{x}}^{(j)} \dots \right. \right. \\
 &\quad \left. \left. f \left(\mathbf{W}_{\mathbf{x}}^{(2)} f \left(\mathbf{W}_{\mathbf{x}}^{(1)} \mathbf{x}_T + \mathbf{W}_{\mathbf{h}}^{(1)} \mathbf{h}_{T-1}^{(1)} + \mathbf{b}_{\mathbf{h}}^{(1)} \right) \right) \right. \right. \\
 &\quad \left. \left. + \mathbf{W}_{\mathbf{h}}^{(2)} \mathbf{h}_{T-1}^{(2)} + \mathbf{b}_{\mathbf{h}}^{(2)} \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{h}}^{(j)} \right) \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \dots f \left(\mathbf{W}_{\mathbf{x}}^{(j)} \dots \right. \right. \\
 &\quad \left. \left. f \left(\mathbf{W}_{\mathbf{x}}^{(2)} \mathbf{h}_T^{(1)} + \mathbf{W}_{\mathbf{h}}^{(2)} \mathbf{h}_{T-1}^{(2)} + \mathbf{b}_{\mathbf{h}}^{(2)} \right) \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{h}}^{(j)} \right) \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= \dots \\
 &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \dots \right. \\
 &\quad \left. f \left(\mathbf{W}_{\mathbf{x}}^{(j)} \mathbf{h}_T^{(j-1)} + \mathbf{W}_{\mathbf{h}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{h}}^{(j)} \right) \right. \\
 &\quad \left. \dots + \mathbf{W}_{\mathbf{h}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_{\mathbf{h}}^{(i)} \right)
 \end{aligned}$$

$$\begin{aligned}
 &= \dots \\
 &= f \left(\mathbf{W}_{\mathbf{x}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{h}}^{(i)} \right) \\
 &= \mathbf{h}_T^{(i)},
 \end{aligned}$$

where we used the definition of the hidden states in Equations (3)-(4). In particular, we have for $i = k$ that $\hat{\mathbf{h}}_{T+k-1}^{\{k\}} = \mathbf{h}_T^{(k)}$.

Now, we show that $\hat{\mathbf{y}}_{T+k-1} = \mathbf{y}_T$. By the definition of the output weights and biases in Equation (8), and by the fact that $\hat{\mathbf{h}}_{T+k-1}^{\{k\}} = \mathbf{h}_T^{(k)}$, we obtain

$$\begin{aligned}
 \hat{\mathbf{y}}_{T+k-1} &= g \left(\hat{\mathbf{W}}_{\mathbf{o}} \hat{\mathbf{h}}_{T+k-1} + \hat{\mathbf{b}}_{\mathbf{o}} \right) \\
 &= g \left(\mathbf{W}_{\mathbf{o}} \hat{\mathbf{h}}_{T+k-1}^{\{k\}} + \mathbf{b}_{\mathbf{o}} \right) \\
 &= g \left(\mathbf{W}_{\mathbf{o}} \mathbf{h}_T^{(k)} + \mathbf{b}_{\mathbf{o}} \right) \\
 &= \mathbf{y}_T,
 \end{aligned} \tag{A.2}$$

which completes the proof. \blacksquare

B. Extension to d-LSTMs

A Long Short-Term Memory recurrent cell (Hochreiter & Schmidhuber, 1997) is given by the introduction of a cell state and a series of gates that control the updates of the states. The cell state together with the gates aim to solve the vanishing gradients problems in the RNN. The LSTM cell is highly popular and we refer to the following implementation:

$$\hat{\mathbf{e}}_t = \sigma \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{e}} \mathbf{x}_t + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{e}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\mathbf{e}} \right), \tag{B.3}$$

$$\hat{\mathbf{f}}_t = \sigma \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{f}} \mathbf{x}_t + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{f}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\mathbf{f}} \right), \tag{B.4}$$

$$\hat{\mathbf{o}}_t = \sigma \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{o}} \mathbf{x}_t + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{o}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\mathbf{o}} \right), \tag{B.5}$$

$$\hat{\mathbf{g}}_t = \tanh \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{c}} \mathbf{x}_t + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{c}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\mathbf{c}} \right), \tag{B.6}$$

$$\hat{\mathbf{c}}_t = \hat{\mathbf{f}}_t \odot \hat{\mathbf{c}}_{t-1} + \hat{\mathbf{e}}_t \odot \hat{\mathbf{g}}_t, \tag{B.7}$$

$$\hat{\mathbf{h}}_t = \hat{\mathbf{o}}_t \odot \tanh(\hat{\mathbf{c}}_t), \tag{B.8}$$

where $\hat{\mathbf{e}}_t$ is the input gate, $\hat{\mathbf{f}}_t$ the forget gate, $\hat{\mathbf{o}}_t$ the output gate, $\hat{\mathbf{g}}_t$ the cell gate, $\hat{\mathbf{c}}_t$ the cell state, and $\hat{\mathbf{h}}_t$ the hidden state. The weight matrices are symbolized $\hat{\mathbf{W}}_{\mathbf{x}\mathbf{a}}$ and $\hat{\mathbf{W}}_{\mathbf{h}\mathbf{a}}$ as well as the bias $\hat{\mathbf{b}}_{\mathbf{a}}$, with $\mathbf{a} \in \{\mathbf{e}, \mathbf{c}, \mathbf{f}, \mathbf{o}\}$ being the respective gate. The symbol \odot represents an element-wise product and $\sigma(\cdot)$ is the sigmoid function.

First, we note that the set of Equations (B.3)-(B.8) can be

expanded into the following two equations:

$$\begin{aligned} \hat{\mathbf{c}}_t = & \sigma \left(\hat{\mathbf{W}}_{\text{xf}} \mathbf{x}_t + \hat{\mathbf{W}}_{\text{hf}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\text{f}} \right) \odot \hat{\mathbf{c}}_{t-1} \\ & + \sigma \left(\hat{\mathbf{W}}_{\text{xe}} \mathbf{x}_t + \hat{\mathbf{W}}_{\text{he}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\text{e}} \right) \\ & \odot \tanh \left(\hat{\mathbf{W}}_{\text{xc}} \mathbf{x}_t + \hat{\mathbf{W}}_{\text{hc}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\text{c}} \right), \end{aligned} \quad (\text{B.9})$$

$$\hat{\mathbf{h}}_t = \sigma \left(\hat{\mathbf{W}}_{\text{xo}} \mathbf{x}_t + \hat{\mathbf{W}}_{\text{ho}} \hat{\mathbf{h}}_{t-1} + \hat{\mathbf{b}}_{\text{o}} \right) \odot \tanh(\hat{\mathbf{c}}_t). \quad (\text{B.10})$$

Rewriting the LSTM Equations (B.3)-(B.8) in this form, allows to remain with the recurrent equations where both $\hat{\mathbf{h}}_t$ and $\hat{\mathbf{c}}_t$ depend on the previous hidden and cell states, $\hat{\mathbf{h}}_{t-1}$ and $\hat{\mathbf{c}}_{t-1}$, and the current input \mathbf{x}_t .

Next, we describe the weight matrices for the single-layer d-LSTM that matches a stacked-LSTM with k layers. The matrices and biases follow the exact same pattern as the RNN proof, being the same for all gates.

$$\hat{\mathbf{W}}_{\text{ha}} = \begin{bmatrix} \mathbf{W}_{\text{ha}}^{(1)} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{W}_{\text{xa}}^{(2)} & \mathbf{W}_{\text{ha}}^{(2)} & & \\ \mathbf{0} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \mathbf{W}_{\text{xa}}^{(i)} & \mathbf{W}_{\text{ha}}^{(i)} & \ddots \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{W}_{\text{xa}}^{(k)} & \mathbf{W}_{\text{ha}}^{(k)} \end{bmatrix} \quad (\text{B.11})$$

$$\hat{\mathbf{b}}_{\text{ha}} = \begin{bmatrix} \mathbf{b}_{\text{ha}}^{(1)} \\ \vdots \\ \mathbf{b}_{\text{ha}}^{(k)} \end{bmatrix}, \quad \hat{\mathbf{W}}_{\text{xa}} = \begin{bmatrix} \mathbf{W}_{\text{xa}}^{(1)} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}, \quad (\text{B.12})$$

where $\hat{\mathbf{W}}_{\text{xa}} \in \mathbb{R}^{kn \times q}$ are the input weights, $\hat{\mathbf{W}}_{\text{ha}} \in \mathbb{R}^{kn \times kn}$ the recurrent weights, $\hat{\mathbf{b}}_{\text{ha}} \in \mathbb{R}^{kn}$ the biases, for gate $\mathbf{a} \in \{\mathbf{e}, \mathbf{c}, \mathbf{o}, \mathbf{f}\}$. We follow the same notation for blocks and layers introduced with Theorem 1. We omit the equations for the output element $\hat{\mathbf{y}}_t$ as they are exactly the same as the RNN in Theorem 1, and thus require the same steps for proving that outputs are equal, i.e., $\hat{\mathbf{y}}_{T+k-1} = \mathbf{y}_T$. Therefore, for the LSTM theorem we will focus on the hidden and cell states.

Theorem 2. *Given an input sequence $\{\mathbf{x}_t\}_{t=1 \dots T}$ and a stacked LSTM with k layers, and initial states $\{\mathbf{h}_0^{(i)}, \mathbf{c}_0^{(i)}\}_{i=1 \dots k}$, the d-LSTM with delay $d = k - 1$, defined by Equations (B.11)-(B.12) and initialized with $\hat{\mathbf{h}}_0$ such that $\hat{\mathbf{h}}_{i-1}^{(i)} = \mathbf{h}_0^{(i)}$, $\forall i = 1 \dots k$ and $\hat{\mathbf{c}}_0$ such that $\hat{\mathbf{c}}_{i-1}^{(i)} = \mathbf{c}_0^{(i)}$, $\forall i = 1 \dots k$, produces the same output sequence but delayed by $k - 1$ timesteps, i.e., $\hat{\mathbf{y}}_{t+k-1} = \mathbf{y}_t$ for all $t = 1 \dots T$. Further, the sequence of hidden and cell states at each layer i are equivalent with delay $i - 1$, i.e.,*

$$\hat{\mathbf{h}}_{t+i-1}^{(i)} = \mathbf{h}_t^{(i)} \text{ and } \hat{\mathbf{c}}_{t+i-1}^{(i)} = \mathbf{c}_t^{(i)} \text{ for all } 1 \leq i \leq k \text{ and } t \geq 1.$$

Proof. We prove Theorem 2 by induction on the sequence length t . First, we show that for $t = 1$ the stacked LSTM and the d-LSTM with the constrained weights are equivalent. Namely, for $t = 1$ we show that the outputs, hidden states and cell states are the same, i.e. $\hat{\mathbf{y}}_k = \mathbf{y}_1$, $\hat{\mathbf{h}}_i^{(i)} = \mathbf{h}_1^{(i)}$, and $\hat{\mathbf{c}}_i^{(i)} = \mathbf{c}_1^{(i)}$, respectively. Without loss of generality, we have for any j in $1 \dots k$ the following:

$$\begin{aligned} \hat{\mathbf{h}}_i^{(j)} &= \sigma \left(\hat{\mathbf{W}}_{\text{xo}}^{(j)} \mathbf{x}_i + \hat{\mathbf{W}}_{\text{ho}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j)} + \hat{\mathbf{b}}_{\text{o}}^{(j)} \right) \odot \tanh \left(\hat{\mathbf{c}}_i^{(j)} \right) \\ &= \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j-1)} + \mathbf{W}_{\text{ho}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j)} + \mathbf{b}_{\text{o}}^{(j)} \right) \\ &\quad \odot \tanh \left(\hat{\mathbf{c}}_i^{(j)} \right) \\ &= \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j-1)} + \mathbf{W}_{\text{ho}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{o}}^{(j)} \right) \\ &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j-1)} + \mathbf{W}_{\text{hf}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{f}}^{(j)} \right) \right) \\ &\quad \odot \hat{\mathbf{c}}_{i-1}^{(j)} \\ &\quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j-1)} + \mathbf{W}_{\text{he}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{e}}^{(j)} \right) \\ &\quad \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(j)} \hat{\mathbf{h}}_{i-1}^{(j-1)} + \mathbf{W}_{\text{hc}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{c}}^{(j)} \right) \\ &= \dots \\ &= \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} \dots \left\{ \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} [\dots \right. \right. \right. \\ &\quad \sigma \left(\mathbf{W}_{\text{xo}}^{(2)} \sigma \left(\mathbf{W}_{\text{xo}}^{(1)} \mathbf{x}_1 + \mathbf{W}_{\text{ho}}^{(1)} \mathbf{h}_0^{(1)} + \mathbf{b}_{\text{o}}^{(1)} \right) \right. \\ &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(1)} \mathbf{x}_1 + \mathbf{W}_{\text{hf}}^{(1)} \mathbf{h}_0^{(1)} + \mathbf{b}_{\text{f}}^{(1)} \right) \right. \\ &\quad \odot \mathbf{c}_0^{(1)} \\ &\quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(1)} \mathbf{x}_1 + \mathbf{W}_{\text{he}}^{(1)} \mathbf{h}_0^{(1)} + \mathbf{b}_{\text{e}}^{(1)} \right) \\ &\quad \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(1)} \mathbf{x}_1 + \mathbf{W}_{\text{hc}}^{(1)} \mathbf{h}_0^{(1)} + \mathbf{b}_{\text{c}}^{(1)} \right) \\ &\quad + \mathbf{W}_{\text{ho}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_{\text{o}}^{(2)} \left. \right) \\ &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(2)} (\dots) + \mathbf{W}_{\text{hf}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_{\text{f}}^{(2)} \right) \right. \\ &\quad \odot \hat{\mathbf{c}}_1^{(2)} \\ &\quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(2)} (\dots) + \mathbf{W}_{\text{he}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_{\text{e}}^{(2)} \right) \\ &\quad \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(2)} (\dots) + \mathbf{W}_{\text{hc}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_{\text{c}}^{(2)} \right) \\ &\quad \dots \left. \right] + \mathbf{W}_{\text{ho}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{o}}^{(j)} \left. \right) \\ &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(j)} [\dots] + \mathbf{W}_{\text{hf}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{f}}^{(j)} \right) \right. \\ &\quad \odot \hat{\mathbf{c}}_{j-1}^{(j)} \\ &\quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(j)} [\dots] + \mathbf{W}_{\text{he}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_{\text{e}}^{(j)} \right) \end{aligned}$$

$$\begin{aligned}
 & \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(j)} [\dots] + \mathbf{W}_{\text{hc}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_c^{(j)} \right) \Big\} \\
 & \dots + \mathbf{W}_{\text{ho}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_o^{(i)} \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(i)} (\dots) + \mathbf{W}_{\text{hf}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_f^{(i)} \right) \right. \\
 & \odot \hat{\mathbf{c}}_{i-1}^{\{i\}} \\
 & \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} (\dots) + \mathbf{W}_{\text{he}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_e^{(i)} \right) \right. \\
 & \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} (\dots) + \mathbf{W}_{\text{hc}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_c^{(i)} \right) \Big\} \\
 = & \sigma \left(\mathbf{W}_{\text{xo}}^{(i)} \dots \left\{ \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} [\dots \right. \right. \right. \\
 & \left. \left. \left. \sigma \left(\mathbf{W}_{\text{xo}}^{(2)} \mathbf{h}_1^{(1)} + \mathbf{W}_{\text{ho}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_o^{(2)} \right) \right. \right. \right. \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(2)} \mathbf{h}_1^{(1)} + \mathbf{W}_{\text{hf}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_f^{(2)} \right) \odot \mathbf{c}_0^{(2)} \right. \\
 & \left. \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(2)} \mathbf{h}_1^{(1)} + \mathbf{W}_{\text{he}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_e^{(2)} \right) \right. \right. \\
 & \left. \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(2)} \mathbf{h}_1^{(1)} + \mathbf{W}_{\text{hc}}^{(2)} \mathbf{h}_0^{(2)} + \mathbf{b}_c^{(2)} \right) \right) \right. \\
 & \left. \dots \right] + \mathbf{W}_{\text{ho}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_o^{(j)} \Big\} \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(j)} [\dots] + \mathbf{W}_{\text{hf}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_f^{(j)} \right) \odot \mathbf{c}_0^{(j)} \right. \\
 & \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(j)} [\dots] + \mathbf{W}_{\text{he}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_e^{(j)} \right) \right. \\
 & \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(j)} [\dots] + \mathbf{W}_{\text{hc}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_c^{(j)} \right) \right) \Big\} \\
 & \dots + \mathbf{W}_{\text{ho}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_o^{(i)} \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(i)} (\dots) + \mathbf{W}_{\text{hf}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_f^{(i)} \right) \odot \mathbf{c}_0^{(i)} \right. \\
 & \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} (\dots) + \mathbf{W}_{\text{he}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_e^{(i)} \right) \right. \\
 & \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} (\dots) + \mathbf{W}_{\text{hc}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_c^{(i)} \right) \right) \Big\} \\
 = & \dots \\
 = & \sigma \left(\mathbf{W}_{\text{xo}}^{(i)} \dots \left\{ \sigma \left(\mathbf{W}_{\text{xo}}^{(j)} \mathbf{h}_1^{(j-1)} + \mathbf{W}_{\text{ho}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_o^{(j)} \right) \right. \right. \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(j)} \mathbf{h}_1^{(j-1)} + \mathbf{W}_{\text{hf}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_f^{(j)} \right) \odot \mathbf{c}_0^{(j)} \right. \\
 & \left. \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(j)} \mathbf{h}_1^{(j-1)} + \mathbf{W}_{\text{he}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_e^{(j)} \right) \right. \right. \\
 & \left. \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(j)} \mathbf{h}_1^{(j-1)} + \mathbf{W}_{\text{hc}}^{(j)} \mathbf{h}_0^{(j)} + \mathbf{b}_c^{(j)} \right) \right) \right. \\
 & \left. \dots + \mathbf{W}_{\text{ho}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_o^{(i)} \right. \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(i)} (\dots) + \mathbf{W}_{\text{hf}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_f^{(i)} \right) \odot \mathbf{c}_0^{(i)} \right. \\
 & \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} (\dots) + \mathbf{W}_{\text{he}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_e^{(i)} \right) \right. \\
 & \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} (\dots) + \mathbf{W}_{\text{hc}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_c^{(i)} \right) \right) \Big\} \\
 = & \dots
 \end{aligned}$$

$$\begin{aligned}
 & = \sigma \left(\mathbf{W}_{\text{xo}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{ho}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_o^{(i)} \right) \\
 & \odot \tanh \left(\sigma \left(\mathbf{W}_{\text{xf}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{hf}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_f^{(i)} \right) \odot \mathbf{c}_0^{(i)} \right. \\
 & \left. + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{he}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_e^{(i)} \right) \right. \\
 & \left. \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{hc}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_c^{(i)} \right) \right) \\
 & = \mathbf{h}_1^{(i)},
 \end{aligned}$$

where we used the initialization assumptions $\hat{\mathbf{h}}_{i-1}^{\{i\}} = \mathbf{h}_0^{(i)}$ and $\hat{\mathbf{c}}_{i-1}^{\{i\}} = \mathbf{c}_0^{(i)}$ for all $i = 1 \dots k$, and the definition of the hidden and cell state in Equations (B.9) and (B.10) for $j-1$ blocks, in the previous steps. In particular, we have for layer k that $\hat{\mathbf{h}}_i^{\{k\}} = \mathbf{h}_1^{(k)}$, and using the same transformations as in (A.1) with RNNs, we obtain $\hat{\mathbf{y}}_k = \mathbf{y}_1$. Furthermore, we obtained that:

$$\begin{aligned}
 \hat{\mathbf{c}}_i^{\{i\}} & = \sigma \left(\hat{\mathbf{W}}_{\text{xf}}^{\{i\}} \mathbf{x}_i + \hat{\mathbf{W}}_{\text{hf}}^{\{i\}} \hat{\mathbf{h}}_{i-1} + \hat{\mathbf{b}}_f^{\{i\}} \right) \odot \hat{\mathbf{c}}_{i-1}^{\{i\}} \\
 & \quad + \sigma \left(\hat{\mathbf{W}}_{\text{xe}}^{\{i\}} \mathbf{x}_i + \hat{\mathbf{W}}_{\text{he}}^{\{i\}} \hat{\mathbf{h}}_{i-1} + \hat{\mathbf{b}}_e^{\{i\}} \right) \\
 & \quad \odot \tanh \left(\hat{\mathbf{W}}_{\text{xc}}^{\{i\}} \mathbf{x}_i + \hat{\mathbf{W}}_{\text{hc}}^{\{i\}} \hat{\mathbf{h}}_{i-1} + \hat{\mathbf{b}}_c^{\{i\}} \right) \\
 & = \sigma \left(\mathbf{W}_{\text{xf}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i-1\}} + \mathbf{W}_{\text{hf}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i\}} + \mathbf{b}_f^{(i)} \right) \odot \mathbf{c}_0^{(i)} \\
 & \quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i-1\}} + \mathbf{W}_{\text{he}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i\}} + \mathbf{b}_e^{(i)} \right) \\
 & \quad \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i-1\}} + \mathbf{W}_{\text{hc}}^{(i)} \hat{\mathbf{h}}_{i-1}^{\{i\}} + \mathbf{b}_c^{(i)} \right) \\
 & = \sigma \left(\mathbf{W}_{\text{xf}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{hf}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_f^{(i)} \right) \odot \mathbf{c}_0^{(i)} \\
 & \quad + \sigma \left(\mathbf{W}_{\text{xe}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{he}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_e^{(i)} \right) \\
 & \quad \odot \tanh \left(\mathbf{W}_{\text{xc}}^{(i)} \mathbf{h}_1^{(i-1)} + \mathbf{W}_{\text{hc}}^{(i)} \mathbf{h}_0^{(i)} + \mathbf{b}_c^{(i)} \right) \\
 & = \mathbf{c}_1^{(i)}
 \end{aligned}$$

Which concludes the basis of the induction.

Next, we assume that $\hat{\mathbf{h}}_{t+i-1}^{\{i\}} = \mathbf{h}_t^{(i)}$ and $\hat{\mathbf{c}}_{t+i-1}^{\{i\}} = \mathbf{c}_t^{(i)}$ for all $1 \leq i \leq k$ and $t \leq T-1$, and prove that it holds for the hidden and cell states for all layers when $t = T$: $\hat{\mathbf{h}}_{T+i-1}^{\{i\}} = \mathbf{h}_T^{(i)}$, $\forall 1 \leq i \leq k$. Without loss of generality, we have for the hidden state $\hat{\mathbf{h}}_{T+i-1}^{\{i\}}$ in constrained weights

$$\begin{aligned}
 &= \dots \\
 &= \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{o}}^{(i)} \dots \left\{ \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{o}}^{(j)} \mathbf{h}_T^{(j-1)} + \mathbf{W}_{\mathbf{h}\mathbf{o}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{o}}^{(j)} \right) \right. \right. \\
 &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{f}}^{(j)} \mathbf{h}_T^{(j-1)} + \mathbf{W}_{\mathbf{h}\mathbf{f}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{f}}^{(j)} \right) \right) \\
 &\quad \odot \mathbf{c}_{T-1}^{(j)} + \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{e}}^{(j)} \mathbf{h}_T^{(j-1)} + \mathbf{W}_{\mathbf{h}\mathbf{e}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{e}}^{(j)} \right) \\
 &\quad \left. \odot \tanh \left(\mathbf{W}_{\mathbf{x}\mathbf{c}}^{(j)} \mathbf{h}_T^{(j-1)} + \mathbf{W}_{\mathbf{h}\mathbf{c}}^{(j)} \mathbf{h}_{T-1}^{(j)} + \mathbf{b}_{\mathbf{c}}^{(j)} \right) \right\} \\
 &\quad \dots + \mathbf{W}_{\mathbf{h}\mathbf{o}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{o}}^{(i)} \\
 &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{f}}^{(i)} (\dots) + \mathbf{W}_{\mathbf{h}\mathbf{f}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{f}}^{(i)} \right) \right) \\
 &\quad \odot \mathbf{c}_{T-1}^{(i)} + \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{e}}^{(i)} (\dots) + \mathbf{W}_{\mathbf{h}\mathbf{e}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{e}}^{(i)} \right) \\
 &\quad \odot \tanh \left(\mathbf{W}_{\mathbf{x}\mathbf{c}}^{(i)} (\dots) + \mathbf{W}_{\mathbf{h}\mathbf{c}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{c}}^{(i)} \right) \\
 &= \dots \\
 &= \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{o}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{o}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{o}}^{(i)} \right) \\
 &\quad \odot \tanh \left(\sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{f}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{f}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{f}}^{(i)} \right) \right) \\
 &\quad \odot \mathbf{c}_{T-1}^{(i)} + \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{e}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{e}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{e}}^{(i)} \right) \\
 &\quad \odot \tanh \left(\mathbf{W}_{\mathbf{x}\mathbf{c}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{c}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{c}}^{(i)} \right) \\
 &= \mathbf{h}_T^{(i)},
 \end{aligned}$$

where we use the recurrent definition of the hidden and cell states in Equations (B.9) and (B.10). In particular, we obtained for $i = k$ that $\hat{\mathbf{h}}_{T+k-1}^{\{k\}} = \mathbf{h}_T^{(k)}$. Applying the same steps as in the d-RNN proof in Eq. (A.2), we obtain $\hat{\mathbf{y}}_{T+k-1} = \mathbf{y}_T$. Last, we obtain for the cell state that

$$\begin{aligned}
 \hat{\mathbf{c}}_{T+i-1}^{\{i\}} &= \sigma \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{f}}^{\{i\}} \mathbf{x}_{T+i-1} + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{f}}^{\{i\}} \hat{\mathbf{h}}_{T+i-2} + \hat{\mathbf{b}}_{\mathbf{f}}^{\{i\}} \right) \\
 &\quad \odot \hat{\mathbf{c}}_{T+i-2}^{\{i\}} \\
 &\quad + \sigma \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{e}}^{\{i\}} \mathbf{x}_{T+i-1} + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{e}}^{\{i\}} \hat{\mathbf{h}}_{T+i-2} + \hat{\mathbf{b}}_{\mathbf{e}}^{\{i\}} \right) \\
 &\quad \odot \tanh \left(\hat{\mathbf{W}}_{\mathbf{x}\mathbf{c}}^{\{i\}} \mathbf{x}_{T+i-1} + \hat{\mathbf{W}}_{\mathbf{h}\mathbf{c}}^{\{i\}} \hat{\mathbf{h}}_{T+i-2} + \hat{\mathbf{b}}_{\mathbf{c}}^{\{i\}} \right) \\
 &= \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{f}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i-1\}} + \mathbf{W}_{\mathbf{h}\mathbf{f}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{f}}^{(i)} \right) \odot \mathbf{c}_{T-1}^{(i)} \\
 &\quad + \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{e}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i-1\}} + \mathbf{W}_{\mathbf{h}\mathbf{e}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{e}}^{(i)} \right) \\
 &\quad \odot \tanh \left(\mathbf{W}_{\mathbf{x}\mathbf{c}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i-1\}} + \mathbf{W}_{\mathbf{h}\mathbf{c}}^{(i)} \hat{\mathbf{h}}_{T+i-2}^{\{i\}} + \mathbf{b}_{\mathbf{c}}^{(i)} \right) \\
 &= \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{f}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{f}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{f}}^{(i)} \right) \odot \mathbf{c}_{T-1}^{(i)} \\
 &\quad + \sigma \left(\mathbf{W}_{\mathbf{x}\mathbf{e}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{e}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{e}}^{(i)} \right) \\
 &\quad \odot \tanh \left(\mathbf{W}_{\mathbf{x}\mathbf{c}}^{(i)} \mathbf{h}_T^{(i-1)} + \mathbf{W}_{\mathbf{h}\mathbf{c}}^{(i)} \mathbf{h}_{T-1}^{(i)} + \mathbf{b}_{\mathbf{c}}^{(i)} \right) \\
 &= \mathbf{c}_T^{(i)}
 \end{aligned}$$

Which completes the proof. \blacksquare

C. Weight Constraints and Connections in d-RNN

Figure 6 shows the weight constraints imposed to achieve equivalence between the stacked RNN and single-layer d-RNN, and a visualization of the d-RNN as connections in the stacked RNN. Figure 6(b) depicts the delay (or ‘‘shift’’) of all the hidden states as they would be computed in the stacked RNN. Each layer is equivalent to a shift by one timestep.

D. Additional Plots for Error Maps

Figure 7 present the standard deviation diagrams for the error maps in Figure 5.

E. Masked Character-Level Language Modeling: Additional Results

In Table 3, we include additional results for smaller networks of the masked language model task. We sampled more delay values for d-LSTMs, but the general conclusions remain the same: intermediate values of delay achieve the lowest BPC. Forward-pass runtimes across delay values show a small increase with larger delays, but the increment is relatively flat compared to stacked LSTMs or (stacked) Bi-LSTMs as they increase in depth. For these experiments, we also used a batch of 128 sequences, and an embedding of dimension 10.

F. Part-of-Speech Tagging: Additional Details and Results

In this section, we include more details about the dataset and the results of all the combinations for the Parts-Of-Speech experiment. We used treebanks from Universal Dependencies (UD) (Nivre et al., 2016) version 2.3. We selected the English EWT treebank¹ (Silveira et al., 2014) (254,854 words), French GSD treebank² (411,465 words), and German GSD treebank³ (297,836 words) based on the quality assigned by the UD authors. We follow the partitioning onto training, validation and test datasets as pre-defined in UD. All treebanks use the same POS tag set containing 17 tags. We use the Polyglot project (Al-Rfou’ et al., 2013) word embeddings (64 dimensions). We build our own alphabets

¹https://github.com/UniversalDependencies/UD_English-EWT/tree/r2.3

²https://github.com/UniversalDependencies/UD_French-GSD/tree/r2.3

³https://github.com/UniversalDependencies/UD_German-GSD/tree/r2.3

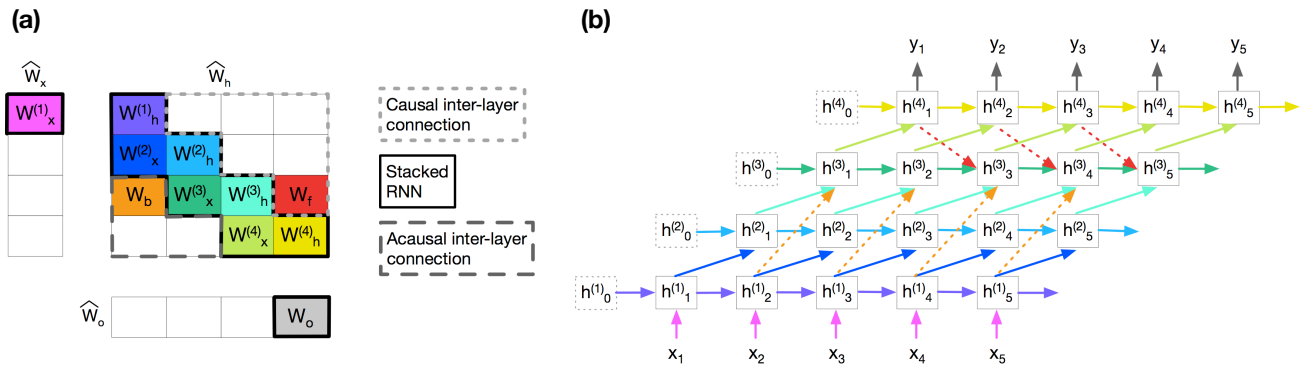


Figure 6. (a) Weights of the single-layer and weight constrained d-RNN that are equivalent to connections in the stacked RNN from Figure 2. (b) Connections in the d-RNN based on the weight matrix in (a). The d-RNN is depicted as it would be the stacked RNN. The hidden states are delayed in time with respect to the stacked network.

Table 3. Performance for smaller networks on the masked character-level language modeling task. Mean and standard deviation values are computed over 5 repetitions of training and inference runtime on the test set.

MODEL	LAYERS	DELAY	UNITS	PARAMS.	VAL. BPC	TEST BPC	RUNTIME
LSTM	1	-	512	1087283	2.139 ± 0.005	2.195 ± 0.002	2.85ms ± 0.14
LSTM	2	-	298	1090689	2.156 ± 0.003	2.215 ± 0.002	6.69ms ± 0.27
LSTM	5	-	172	1083735	2.199 ± 0.016	2.255 ± 0.015	11.32ms ± 0.05
Bi-LSTM	1	-	360	1091107	1.130 ± 0.003	1.187 ± 0.004	5.82ms ± 0.18
Bi-LSTM	2	-	182	1090487	0.800 ± 0.004	0.846 ± 0.005	11.08ms ± 0.59
Bi-LSTM	5	-	102	1104151	0.796 ± 0.007	0.841 ± 0.006	23.94ms ± 0.17
D-LSTM	1	1	512	1087283	1.470 ± 0.002	1.518 ± 0.003	2.80ms ± 0.02
D-LSTM	1	2	512	1087283	1.162 ± 0.004	1.208 ± 0.003	2.81ms ± 0.01
D-LSTM	1	3	512	1087283	0.995 ± 0.002	1.039 ± 0.002	3.02ms ± 0.23
D-LSTM	1	5	512	1087283	0.877 ± 0.001	0.920 ± 0.003	3.01ms ± 0.22
D-LSTM	1	8	512	1087283	0.859 ± 0.002	0.905 ± 0.003	3.04ms ± 0.19
D-LSTM	1	10	512	1087283	0.889 ± 0.004	0.935 ± 0.005	3.22ms ± 0.18
D-LSTM	1	15	512	1087283	0.971 ± 0.004	1.014 ± 0.002	3.17ms ± 0.05

based on the most frequent 100 characters in the vocabularies. All the networks have a 100-dimensional character-level embedding, which is trained with the network. We use a batch size of 32 sentences.

Results for German, English, and French can be found in Tables 4, 5, and 6, respectively. The best result that does not use a bidirectional network is marked in bold for each language.

References

Al-Rfou', R., Perozzi, B., and Skiena, S. Polyglot: Distributed word representations for multilingual NLP. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pp. 183–192, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W13-3520>.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>.

Nivre, J., de Marneffe, M.-C., Ginter, F., Goldberg, Y., Hajič, J., Manning, C. D., McDonald, R., Petrov, S., Pyysalo, S., Silveira, N., Tsarfaty, R., and Zeman, D. Universal dependencies v1: A multilingual treebank collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pp. 1659–1666, Portorož, Slovenia, May 2016. European Language Resources Association (ELRA). URL <https://www.aclweb.org/anthology/L16-1262>.

Silveira, N., Dozat, T., de Marneffe, M.-C., Bowman, S., Connor, M., Bauer, J., and Manning, C. A gold standard dependency corpus for English. In *Proceedings of the Ninth International Conference on*

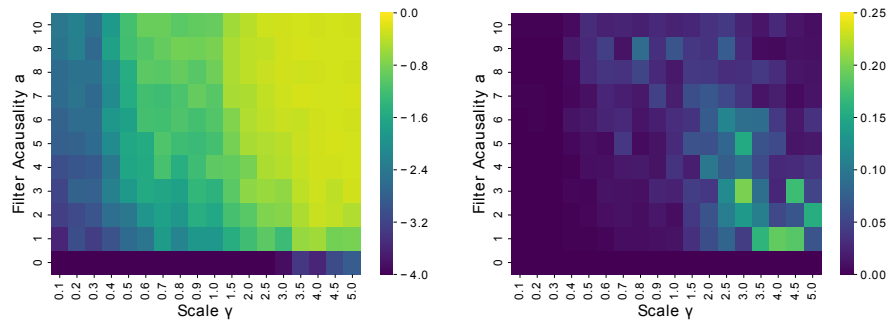
Table 4. Parts-of-Speech results for German. The table shows all possible combinations of delays or bidirectional LSTM networks. The best forward-only network is marked in bold.

CHARACTER-LEVEL NETWORK	WORD-LEVEL NETWORK	VALIDATION ACCURACY	TEST ACCURACY
Bi-LSTM	Bi-LSTM	93.88 ± 0.13	93.15 ± 0.08
Bi-LSTM	LSTM	92.00 ± 0.16	91.50 ± 0.05
Bi-LSTM	D-LSTM WITH DELAY=1	93.32 ± 0.23	92.81 ± 0.14
Bi-LSTM	D-LSTM WITH DELAY=2	93.15 ± 0.06	92.67 ± 0.08
Bi-LSTM	D-LSTM WITH DELAY=3	92.82 ± 0.14	92.25 ± 0.16
Bi-LSTM	D-LSTM WITH DELAY=4	92.41 ± 0.12	91.95 ± 0.17
Bi-LSTM	D-LSTM WITH DELAY=5	91.86 ± 0.11	91.57 ± 0.20
LSTM	Bi-LSTM	93.96 ± 0.12	93.43 ± 0.07
LSTM	LSTM	92.05 ± 0.16	91.58 ± 0.11
LSTM	D-LSTM WITH DELAY=1	93.46 ± 0.16	92.71 ± 0.11
LSTM	D-LSTM WITH DELAY=2	93.13 ± 0.10	92.61 ± 0.26
LSTM	D-LSTM WITH DELAY=3	92.91 ± 0.13	92.38 ± 0.15
LSTM	D-LSTM WITH DELAY=4	92.56 ± 0.17	92.06 ± 0.19
D-LSTM WITH DELAY=1	Bi-LSTM	93.93 ± 0.06	93.39 ± 0.18
D-LSTM WITH DELAY=1	LSTM	92.04 ± 0.11	91.58 ± 0.14
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=1	93.48 ± 0.31	92.87 ± 0.24
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=2	93.11 ± 0.18	92.54 ± 0.08
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=3	92.85 ± 0.14	92.28 ± 0.19
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=4	92.50 ± 0.12	92.11 ± 0.19
D-LSTM WITH DELAY=3	Bi-LSTM	94.00 ± 0.17	93.32 ± 0.18
D-LSTM WITH DELAY=3	LSTM	92.10 ± 0.24	91.61 ± 0.18
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=1	93.29 ± 0.09	92.68 ± 0.09
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=2	93.09 ± 0.21	92.59 ± 0.16
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=3	92.86 ± 0.24	92.42 ± 0.16
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=4	92.53 ± 0.17	92.08 ± 0.18
D-LSTM WITH DELAY=5	Bi-LSTM	93.88 ± 0.17	93.27 ± 0.06
D-LSTM WITH DELAY=5	LSTM	91.88 ± 0.18	91.54 ± 0.11
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=1	93.31 ± 0.14	92.74 ± 0.10
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=2	93.17 ± 0.13	92.57 ± 0.17
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=3	92.84 ± 0.19	92.25 ± 0.10
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=4	92.50 ± 0.22	91.96 ± 0.19

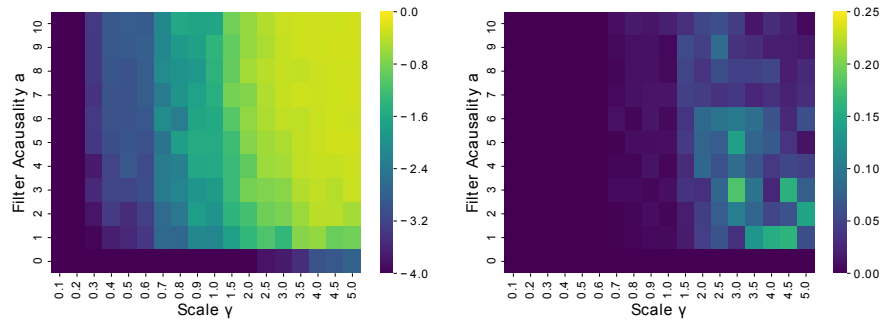
Language Resources and Evaluation (LREC'14), pp. 2897–2904, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2014/pdf/1089_Paper.pdf.

Table 5. Parts-of-Speech results for English. The table shows all possible combinations of delays or bidirectional LSTM networks. The best forward-only network is marked in bold.

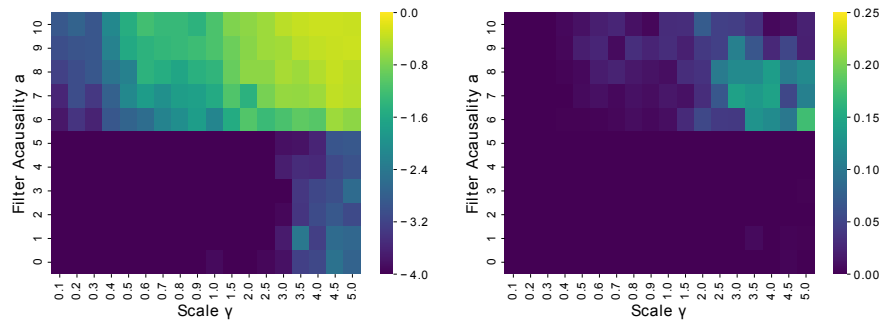
CHARACTER-LEVEL NETWORK	WORD-LEVEL NETWORK	VALIDATION ACCURACY	TEST ACCURACY
Bi-LSTM	Bi-LSTM	94.85 ± 0.05	94.84 ± 0.08
Bi-LSTM	LSTM	91.90 ± 0.12	92.05 ± 0.09
Bi-LSTM	D-LSTM WITH DELAY=1	94.47 ± 0.06	94.41 ± 0.05
Bi-LSTM	D-LSTM WITH DELAY=2	94.17 ± 0.13	94.14 ± 0.10
Bi-LSTM	D-LSTM WITH DELAY=3	93.70 ± 0.07	93.87 ± 0.07
Bi-LSTM	D-LSTM WITH DELAY=4	93.11 ± 0.14	93.26 ± 0.08
Bi-LSTM	D-LSTM WITH DELAY=5	92.54 ± 0.16	92.70 ± 0.10
LSTM	Bi-LSTM	95.03 ± 0.14	94.99 ± 0.15
LSTM	LSTM	92.05 ± 0.13	92.14 ± 0.10
LSTM	D-LSTM WITH DELAY=1	94.53 ± 0.08	94.58 ± 0.11
LSTM	D-LSTM WITH DELAY=2	94.29 ± 0.05	94.28 ± 0.05
LSTM	D-LSTM WITH DELAY=3	93.81 ± 0.11	93.85 ± 0.12
LSTM	D-LSTM WITH DELAY=4	93.39 ± 0.12	93.55 ± 0.10
D-LSTM WITH DELAY=1	Bi-LSTM	94.94 ± 0.07	94.95 ± 0.06
D-LSTM WITH DELAY=1	LSTM	91.96 ± 0.16	92.09 ± 0.10
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=1	94.57 ± 0.08	94.57 ± 0.14
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=2	94.29 ± 0.12	94.37 ± 0.08
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=3	93.86 ± 0.05	93.84 ± 0.10
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=4	93.35 ± 0.10	93.56 ± 0.13
D-LSTM WITH DELAY=3	Bi-LSTM	94.98 ± 0.09	94.91 ± 0.10
D-LSTM WITH DELAY=3	LSTM	91.96 ± 0.08	92.08 ± 0.10
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=1	94.47 ± 0.03	94.51 ± 0.10
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=2	94.21 ± 0.05	94.18 ± 0.03
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=3	93.80 ± 0.13	93.88 ± 0.13
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=4	93.23 ± 0.13	93.38 ± 0.11
D-LSTM WITH DELAY=5	Bi-LSTM	94.90 ± 0.07	94.87 ± 0.09
D-LSTM WITH DELAY=5	LSTM	91.84 ± 0.11	91.98 ± 0.20
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=1	94.36 ± 0.09	94.44 ± 0.08
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=2	94.05 ± 0.07	94.19 ± 0.05
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=3	93.61 ± 0.07	93.76 ± 0.05
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=4	93.14 ± 0.04	93.27 ± 0.12



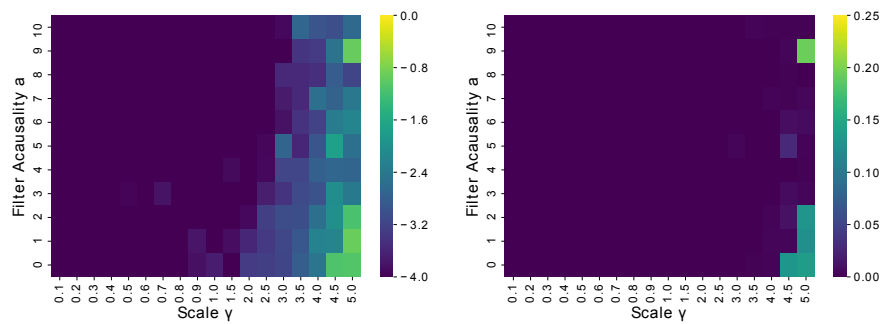
(a) LSTM



(b) Bi-LSTM



(c) d-LSTM with delay=5



(d) d-LSTM with delay=10

Figure 7. Error maps presented in Figure 4 (left column) together with their standard deviation figures.

Table 6. Parts-of-Speech results for French. The table shows all possible combinations of delays or bidirectional LSTM networks. The best forward-only network is marked in bold.

CHARACTER-LEVEL NETWORK	WORD-LEVEL NETWORK	VALIDATION ACCURACY	TEST ACCURACY
Bi-LSTM	Bi-LSTM	97.63 ± 0.06	97.22 ± 0.11
Bi-LSTM	LSTM	96.67 ± 0.05	96.15 ± 0.17
Bi-LSTM	D-LSTM WITH DELAY=1	97.48 ± 0.02	96.98 ± 0.05
Bi-LSTM	D-LSTM WITH DELAY=2	97.41 ± 0.02	96.91 ± 0.12
Bi-LSTM	D-LSTM WITH DELAY=3	97.31 ± 0.05	96.84 ± 0.09
Bi-LSTM	D-LSTM WITH DELAY=4	97.12 ± 0.05	96.61 ± 0.06
Bi-LSTM	D-LSTM WITH DELAY=5	96.88 ± 0.10	96.20 ± 0.14
LSTM	Bi-LSTM	97.70 ± 0.07	97.19 ± 0.09
LSTM	LSTM	96.67 ± 0.07	96.10 ± 0.11
LSTM	D-LSTM WITH DELAY=1	97.49 ± 0.07	97.03 ± 0.07
LSTM	D-LSTM WITH DELAY=2	97.49 ± 0.05	97.00 ± 0.06
LSTM	D-LSTM WITH DELAY=3	97.34 ± 0.04	96.89 ± 0.09
LSTM	D-LSTM WITH DELAY=4	97.16 ± 0.06	96.66 ± 0.15
D-LSTM WITH DELAY=1	Bi-LSTM	97.67 ± 0.07	97.23 ± 0.12
D-LSTM WITH DELAY=1	LSTM	96.66 ± 0.06	95.97 ± 0.07
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=1	97.49 ± 0.04	97.04 ± 0.13
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=2	97.43 ± 0.05	96.98 ± 0.05
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=3	97.36 ± 0.08	96.80 ± 0.10
D-LSTM WITH DELAY=1	D-LSTM WITH DELAY=4	97.22 ± 0.06	96.57 ± 0.10
D-LSTM WITH DELAY=3	Bi-LSTM	97.67 ± 0.08	97.21 ± 0.08
D-LSTM WITH DELAY=3	LSTM	96.67 ± 0.07	95.98 ± 0.14
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=1	97.52 ± 0.04	97.02 ± 0.09
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=2	97.44 ± 0.02	96.97 ± 0.12
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=3	97.28 ± 0.04	96.74 ± 0.07
D-LSTM WITH DELAY=3	D-LSTM WITH DELAY=4	97.13 ± 0.05	96.57 ± 0.09
D-LSTM WITH DELAY=5	Bi-LSTM	97.61 ± 0.03	97.12 ± 0.06
D-LSTM WITH DELAY=5	LSTM	96.64 ± 0.06	96.08 ± 0.08
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=1	97.46 ± 0.02	96.96 ± 0.13
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=2	97.41 ± 0.06	96.87 ± 0.06
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=3	97.36 ± 0.05	96.82 ± 0.07
D-LSTM WITH DELAY=5	D-LSTM WITH DELAY=4	97.15 ± 0.05	96.51 ± 0.07