# Long-short Term Memory (LSTM): back-propagation derivation

Harsha Vardhan

May 20, 2022

**Abstract**

This document contains derivation of the gradients for a Long-short Term Memory Unit, using back-propagation (or reverse-mode differentiation). For the implementation of the neural-network see the accompanying notebooks.
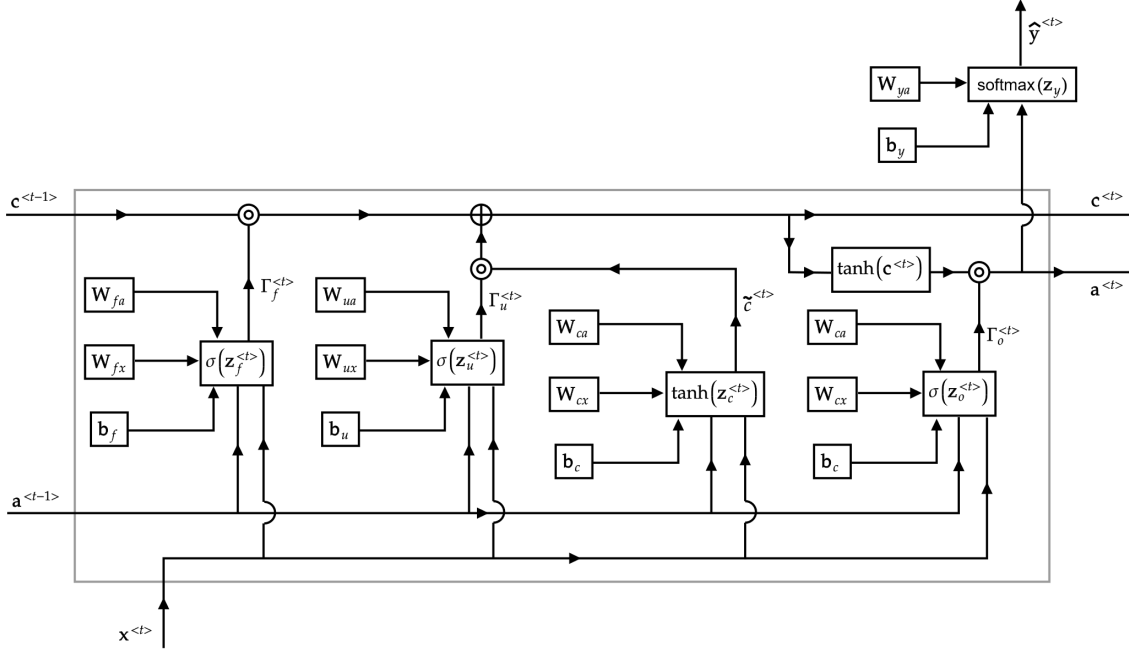
## 1 Block Architecture



Figure 1: forward propagation diagram for a LSTM-block at time-step $t$. The concentric circles represent a hadamard-product (i.e. $\mathbf{x} \circ \mathbf{y}$) of the input vectors.

## 2 Forward Propagation

Given, $\mathbf{c}^{\langle t-1 \rangle}$ and $\mathbf{a}^{\langle t-1 \rangle}$ from the $(t-1)^{th}$ time-step, the equations for forward propagation through the $t^{th}$ time-step, are as follows: (see figure-1)

1

$$\mathbf{z}_f^{\langle t \rangle} = \mathbf{W}_{fa}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{fx}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{b}_f \quad \text{eq.1}$$

$$\Gamma_f^{\langle t \rangle} = \sigma(\mathbf{z}_f^{\langle t \rangle}) \quad \text{eq.2}$$

$$\mathbf{z}_u^{\langle t \rangle} = \mathbf{W}_{ua}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{ux}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{b}_u \quad \text{eq.3}$$

$$\Gamma_u^{\langle t \rangle} = \sigma(\mathbf{z}_u^{\langle t \rangle}) \quad \text{eq.4}$$

$$\mathbf{z}_o^{\langle t \rangle} = \mathbf{W}_{oa}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{ox}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{b}_o \quad \text{eq.5}$$

$$\Gamma_o^{\langle t \rangle} = \sigma(\mathbf{z}_o^{\langle t \rangle}) \quad \text{eq.6}$$

$$\mathbf{z}_c^{\langle t \rangle} = \mathbf{W}_{ca}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{cx}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{b}_c \quad \text{eq.7}$$

$$\tilde{\mathbf{c}}^{\langle t \rangle} = \tanh(\mathbf{z}_c^{\langle t \rangle}) \quad \text{eq.8}$$

$$\mathbf{c}^{\langle t \rangle} = \Gamma_u \circ \tilde{\mathbf{c}}^{\langle t \rangle} + \Gamma_f \circ \mathbf{c}^{\langle t-1 \rangle} \quad \text{eq.9}$$

$$\mathbf{a}^{\langle t \rangle} = \tanh(\mathbf{c}^{\langle t \rangle}) \circ \Gamma_o \quad \text{eq.10}$$

$$\mathbf{z}_y^{\langle t \rangle} = \mathbf{W}_{ya}^\intercal \mathbf{a}^{\langle t \rangle} + \mathbf{b}_y \quad \text{eq.11}$$

$$\hat{\mathbf{y}}^{\langle t \rangle} = \text{softmax}(\mathbf{z}_y^{\langle t \rangle}) \quad \text{eq.12}$$

where,

$\mathbf{x}^{\langle t \rangle}$ is a $(n_x, 1)$-dimensional vector

$\mathbf{c}^{\langle t \rangle}$ & $\mathbf{c}^{\langle t-1 \rangle}$ & $\tilde{\mathbf{c}}^{\langle t \rangle}$ are $(n_a, 1)$-dimensional vectors

$\mathbf{a}^{\langle t \rangle}$ & $\mathbf{a}^{\langle t-1 \rangle}$ are $(n_a, 1)$-dimensional vectors

$\mathbf{W}_{*a}, * \in \{f, u, o, c\}$ are $(n_a, n_a)$-dimensional parameter matrices

$\mathbf{W}_{*x}, * \in \{f, u, o, c\}$ are $(n_x, n_a)$-dimensional parameter matrices

$\mathbf{b}_{*}, * \in \{f, u, o, c\}$ are $(n_a, 1)$-dimensional bias-vectors

$\Gamma_*^{\langle t \rangle}, * \in \{f, u, o\}$ are $(n_a, 1)$-dimensional vector, which represents a gate

$\mathbf{z}_*^{\langle t \rangle}, * \in \{f, u, o, c\}$ are $(n_a, 1)$-dimensional vector, is input to a gates' activation

$\mathbf{z}_y^{\langle t \rangle}$ is a $(n_y, 1)$-dimensional vector

$\mathbf{W}_{ya}$ is a $(n_a, n_y)$-dimensional parameter-matrix

$\mathbf{b}_y$ is a $(n_y, 1)$-dimensional bias vector

Also, the sub-scripts denote the following: $_o$ - output-gate, $_u$ - update-gate, and $_f$ - forget-gate.
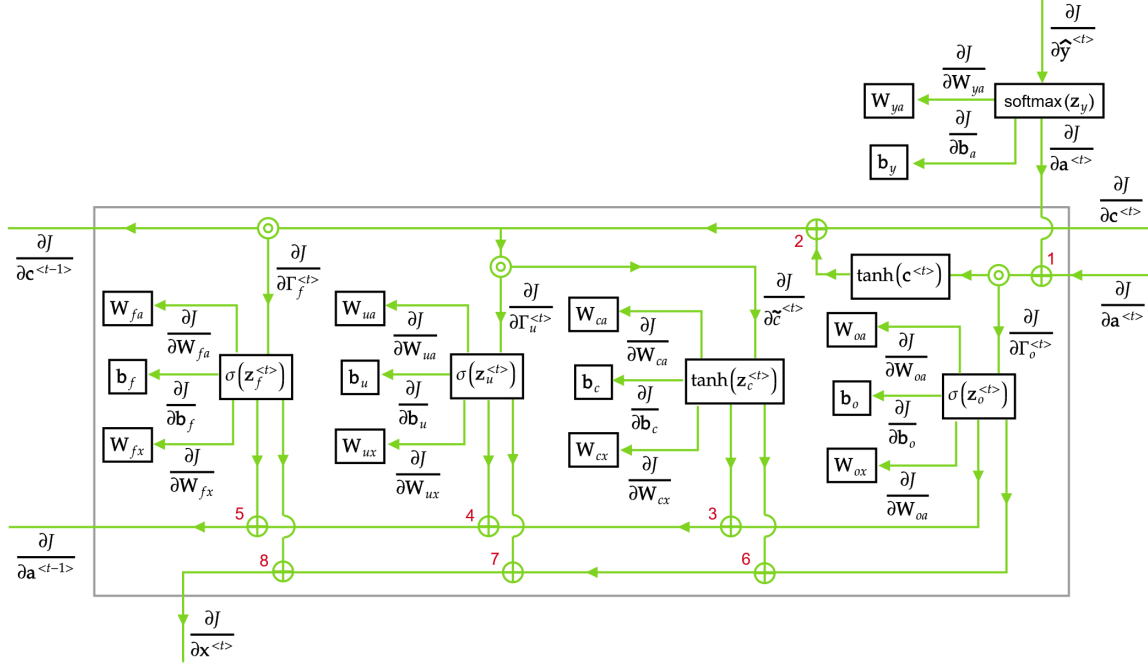
## 2.1  Backward Propagation



Figure 2: backward-propagation/gradient-flow diagram for a LSTM-block at time-step $t$. The concentric circles represent a hadamard-product (i.e. $\mathbf{x} \circ \mathbf{y}$) of the input vectors.

### 2.1.1  Computing $\frac{\partial J}{\partial \hat{\mathbf{y}}^{\langle t \rangle}}$

Since, $\hat{\mathbf{y}}^{\langle t \rangle}$ is computed using the softmax() activation-function, the loss $J$ is computed using the cross-entropy loss. Also, let $\mathbf{y}^{\langle t \rangle}$ be the output-label corresponding to the $t^{th}$ time-step. Then,

$$\frac{\partial J}{\partial \hat{\mathbf{y}}^{\langle t \rangle}} = \begin{bmatrix} \frac{\partial J}{\partial \hat{y}_1^{\langle t \rangle}} & \frac{\partial J}{\partial \hat{y}_2^{\langle t \rangle}} & \cdots & \frac{\partial J}{\partial \hat{y}_{n_y}^{\langle t \rangle}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{y_1^{\langle t \rangle}}{\hat{y}_1^{\langle t \rangle}} & \frac{y_2^{\langle t \rangle}}{\hat{y}_2^{\langle t \rangle}} & \cdots & \frac{y_{n_y}^{\langle t \rangle}}{\hat{y}_{n_y}^{\langle t \rangle}} \end{bmatrix}$$

### 2.1.2  Computing $\frac{\partial J}{\partial \mathbf{W}_{ya}}$ and $\frac{\partial J}{\partial \mathbf{b}_y}$

From *eq.12* (see section-2), we have

$$\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}} = \frac{\partial J}{\partial \hat{\mathbf{y}}^{\langle t \rangle}} \frac{\partial \hat{\mathbf{y}}^{\langle t \rangle}}{\partial \mathbf{z}_y^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \hat{\mathbf{y}}^{\langle t \rangle}} \left( \text{diag}(\hat{\mathbf{y}}^{\langle t \rangle}) - \hat{\mathbf{y}}^{\langle t \rangle}(\hat{\mathbf{y}}^{\langle t \rangle})^{\mathsf{T}} \right)$$

From *eq.11* (see section-2), we have

$$\text{tr}(J) = \text{tr}\left( \frac{\partial J}{\partial \mathbf{z}_y^{\langle r \rangle}} d\mathbf{z}_y^{\langle t \rangle} \right)$$

3

where, $d\mathbf{z}_y^{\langle t \rangle}$ can be expanded as follows:

$$d\mathbf{z}_y^{\langle t \rangle} = d\left(\mathbf{W}_{ya}^{\mathsf{T}}\mathbf{a}^{\langle t \rangle} + \mathbf{b}_y\right)$$
$$= d\left(\mathbf{W}_{ya}^{\mathsf{T}}\right)\mathbf{a}^{\langle t \rangle} + \mathbf{W}_{ya}^{\mathsf{T}}d\left(\mathbf{a}^{\langle t \rangle}\right) + d\mathbf{b}_y \qquad \text{eq.13-1}$$

when differentiating w.r.t. $\mathbf{W}_{ya}$, we have $d\mathbf{a}^{\langle t \rangle} = 0$, and $d\mathbf{b}_y = 0$. So,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}} d\mathbf{W}_{ya}^{\mathsf{T}}\mathbf{a}^{\langle t \rangle}\right)$$

$$= \text{tr}\left(\left(\mathbf{a}^{\langle t \rangle}\right)^{\mathsf{T}} d\mathbf{W}_{ya}\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}\right)^{\mathsf{T}}\right)$$

$$= \text{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}\right)^{\mathsf{T}}\left(\mathbf{a}^{\langle t \rangle}\right)^{\mathsf{T}} d\mathbf{W}_{ya}\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{W}_{ya}} = \left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}\right)^{\mathsf{T}}\left(\mathbf{a}^{\langle t \rangle}\right)^{\mathsf{T}}$$

when differentiating w.r.t. $\mathbf{b}_y$, we have $d\mathbf{a}^{\langle t \rangle} = 0$, and $d\mathbf{W}_{ya} = 0$. So,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}} d\mathbf{b}_y^{\mathsf{T}}\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{b}_y} = \frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}$$

### 2.1.3 Computing $\frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}}$

This derivative has two components

**Comp-1** one flows-in from the $(t+1)^{th}$ time-step, and

**Comp-2** one flows-in as the derivative from $\hat{\mathbf{y}}^{\langle t \rangle}$. This derivative is computed from *eq.13-1* as follows,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}} d\mathbf{z}_y^{\langle t \rangle}\right)$$

$$= \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}\mathbf{W}_{ya}^{\mathsf{T}} d\mathbf{a}^{\langle t \rangle}\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{z}_y^{\langle t \rangle}}\mathbf{W}_{ya}^{\mathsf{T}}$$

Then, these two components are added to compute the true derivative (see the $\oplus$ labeled as *1* in the figure-2), i.e.

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} = \left.\frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}}\right|_{\textbf{Comp-1}} + \left.\frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}}\right|_{\textbf{Comp-2}}$$

### 2.1.4 Computing $\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}}$, and $\frac{\partial J}{\partial \Gamma_o^{\langle t \rangle}}$

From *eq.10* (see section-2), we have

$$
\begin{aligned}
\frac{\partial J}{\partial \Gamma_o^{\langle t \rangle}} &= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \frac{\partial \mathbf{a}^{\langle t \rangle}}{\partial \Gamma_o^{\langle t \rangle}} \\
&= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \operatorname{diag}\left(\tanh(\mathbf{c}^{\langle t \rangle})\right) \\
&= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \circ \left(\tanh(\mathbf{c}^{\langle t \rangle})\right)^{\mathsf{T}}
\end{aligned}
$$

**Note**: For more information about computing the derivative of a Hadamard-product, see section-4.

The derivative $\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}}$ is made of two components, described as follows

**Comp-1** one flows-in from the $(t+1)^{th}$ time-step, and

**Comp-2** one flows-in as the derivative from the $\mathbf{a}^{\langle t \rangle}$. This derivative is computed as follows,

$$
\begin{aligned}
\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} &= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \frac{\partial \mathbf{a}^{\langle t \rangle}}{\partial \mathbf{c}^{\langle t \rangle}} \\
&= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \frac{\partial \mathbf{a}^{\langle t \rangle}}{\partial \tanh(\mathbf{c}^{\langle t \rangle})} \frac{\partial \tanh(\mathbf{c}^{\langle t \rangle})}{\partial \mathbf{c}^{\langle t \rangle}} \\
&= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} (\mathbf{I}_{n_a \times n_a} \circ \Gamma_o^{\langle t \rangle})(\mathbf{I}_{n_a \times n_a} - \mathbf{I}_{n_a \times n_a} \circ \tanh^2(\mathbf{c}^{\langle t \rangle})) \\
&= \frac{\partial J}{\partial \mathbf{a}^{\langle t \rangle}} \circ \left[\Gamma_o^{\langle t \rangle} \circ (\mathbf{1}_{(n_a,1)} - \tanh^2(\mathbf{c}^{\langle t \rangle}))\right]^{\mathsf{T}}
\end{aligned}
$$

These two components are then added to compute the true derivative (see $\oplus$ labeled as *2* in figure-2), i.e.

$$
\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} = \left.\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}}\right|_{\textbf{Comp-1}} + \left.\frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}}\right|_{\textbf{Comp-2}}
$$

### 2.1.5 Computing $\frac{\partial J}{\partial \mathbf{W}_{oa}}$, $\frac{\partial J}{\partial \mathbf{W}_{ox}}$, and $\frac{\partial J}{\partial \mathbf{b}_o}$

From *eq.6* (see section-2), we have

$$
\begin{aligned}
\frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} &= \frac{\partial J}{\partial \Gamma_o^{\langle t \rangle}} \frac{\partial \Gamma_o^{\langle t \rangle}}{\partial \mathbf{z}_o^{\langle t \rangle}} \\
&= \frac{\partial J}{\partial \Gamma_o^{\langle t \rangle}} \circ \left(\Gamma_o^{\langle t \rangle} \circ (\mathbf{1}_{(n_a,1)} - \Gamma_o^{\langle t \rangle})\right)^{\mathsf{T}}
\end{aligned}
$$

From *eq.5* (see section-2), using the trace-method, we have

$$
\operatorname{tr}(J) = \operatorname{tr}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} \, d\mathbf{z}_o^{\langle t \rangle}\right)
$$

where, $\mathrm{d}\mathbf{z}_o^{\langle t\rangle}$ can be expanded as follows:

$$\begin{aligned}\mathrm{d}\mathbf{z}_o^{\langle t\rangle} &= \mathrm{d}\left(\mathbf{W}_{oa}^{\mathsf{T}}\,\mathbf{a}^{\langle t-1\rangle} + \mathbf{W}_{ox}^{\mathsf{T}}\,\mathbf{x}^{\langle t\rangle} + \mathbf{b}_o\right)\\ &= \mathrm{d}\mathbf{W}_{oa}^{\mathsf{T}}\,\mathbf{a}^{\langle t-1\rangle} + \mathbf{W}_{oa}^{\mathsf{T}}\,\mathrm{d}\mathbf{a}^{\langle t-1\rangle} + \mathrm{d}\mathbf{W}_{ox}^{\mathsf{T}}\,\mathbf{x}^{\langle t\rangle} + \mathbf{W}_{ox}^{\mathsf{T}}\,\mathrm{d}\mathbf{x}^{\langle t\rangle} + \mathrm{d}\mathbf{b}_o \qquad \text{eq.13-2}\end{aligned}$$

when differentiating w.r.t. $\mathbf{W}_{oa}$, we have $\mathrm{d}\mathbf{W}_{ox} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1\rangle} = 0$, $\mathrm{d}\mathbf{b}_o = 0$, and $\mathrm{d}\mathbf{x}^{\langle t\rangle} = 0$. So,

$$\begin{aligned}\mathrm{d}J &= \mathrm{tr}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\,\mathrm{d}\mathbf{W}_{oa}^{\mathsf{T}}\,\mathbf{a}^{\langle t-1\rangle}\right)\\ &= \mathrm{tr}\left(\left(\mathbf{a}^{\langle t-1\rangle}\right)^{\mathsf{T}}\,\mathrm{d}\mathbf{W}_{oa}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\right)\\ &= \mathrm{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\left(\mathbf{a}^{\langle t-1\rangle}\right)^{\mathsf{T}}\,\mathrm{d}\mathbf{W}_{oa}\right)\\ \implies \frac{\partial J}{\partial \mathbf{W}_{oa}} &= \left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\left(\mathbf{a}^{\langle t-1\rangle}\right)^{\mathsf{T}}\end{aligned}$$

when differentiating w.r.t. $\mathbf{W}_{ox}$, we have $\mathrm{d}\mathbf{W}_{oa} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1\rangle} = 0$, $\mathrm{d}\mathbf{b}_o = 0$, and $\mathrm{d}\mathbf{x}^{\langle t\rangle} = 0$. So,

$$\begin{aligned}\mathrm{d}J &= \mathrm{tr}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\,\mathrm{d}\mathbf{W}_{ox}^{\mathsf{T}}\,\mathbf{x}^{\langle t\rangle}\right)\\ &= \mathrm{tr}\left(\left(\mathbf{x}^{\langle t\rangle}\right)^{\mathsf{T}}\,\mathrm{d}\mathbf{W}_{ox}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\right)\\ &= \mathrm{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\left(\mathbf{x}^{\langle t\rangle}\right)^{\mathsf{T}}\,\mathrm{d}\mathbf{W}_{ox}\right)\\ \implies \frac{\partial J}{\partial \mathbf{W}_{ox}} &= \left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\right)^{\mathsf{T}}\left(\mathbf{x}^{\langle t\rangle}\right)^{\mathsf{T}}\end{aligned}$$

when differentiating w.r.t. $\mathbf{b}_o$, we have $\mathrm{d}\mathbf{W}_{oa} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1\rangle} = 0$, $\mathrm{d}\mathbf{W}_{ox} = 0$, and $\mathrm{d}\mathbf{x}^{\langle t\rangle} = 0$. So,

$$\begin{aligned}\mathrm{d}J &= \mathrm{tr}\left(\frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\,\mathrm{d}\mathbf{b}_o\right)\\ \implies \frac{\partial J}{\partial \mathbf{b}_o} &= \frac{\partial J}{\partial \mathbf{z}_o^{\langle t\rangle}}\end{aligned}$$

### 2.1.6 Computing $\frac{\partial J}{\partial \Gamma_u^{\langle t\rangle}}$, $\frac{\partial J}{\partial \tilde{\mathbf{c}}^{\langle t\rangle}}$, $\frac{\partial J}{\partial \Gamma_f^{\langle t\rangle}}$, and $\frac{\partial J}{\partial \mathbf{c}^{\langle t-1\rangle}}$

From *eq.9* (see section-2), we have

$$\begin{aligned}\frac{\partial J}{\partial \Gamma_u^{\langle t\rangle}} &= \frac{\partial J}{\partial \mathbf{c}^{\langle t\rangle}}\frac{\partial \mathbf{c}^{\langle t\rangle}}{\partial \Gamma_u^{\langle t\rangle}}\\ &= \frac{\partial J}{\partial \mathbf{c}^{\langle t\rangle}}\,\mathrm{diag}\left(\tilde{\mathbf{c}}^{\langle t\rangle}\right)\\ &= \frac{\partial J}{\partial \mathbf{c}^{\langle t\rangle}}\circ\left(\tilde{\mathbf{c}}^{\langle t\rangle}\right)^{\mathsf{T}}\end{aligned}$$

$$\frac{\partial J}{\partial \tilde{\mathbf{c}}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \frac{\partial \mathbf{c}^{\langle t \rangle}}{\partial \tilde{\mathbf{c}}^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \operatorname{diag}\left(\Gamma_u^{\langle t \rangle}\right)$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \circ \left(\Gamma_u^{\langle t \rangle}\right)^{\mathsf{T}}$$

$$\frac{\partial J}{\partial \Gamma_f^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \frac{\partial \mathbf{c}^{\langle t \rangle}}{\partial \Gamma_f^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \operatorname{diag}\left(\mathbf{c}^{\langle t-1 \rangle}\right)$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \circ \left(\mathbf{c}^{\langle t-1 \rangle}\right)^{\mathsf{T}}$$

$$\frac{\partial J}{\partial \mathbf{c}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \frac{\partial \mathbf{c}^{\langle t \rangle}}{\partial \mathbf{c}^{\langle t-1 \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \operatorname{diag}\left(\Gamma_f^{\langle t \rangle}\right)$$

$$= \frac{\partial J}{\partial \mathbf{c}^{\langle t \rangle}} \circ \left(\Gamma_f^{\langle t \rangle}\right)^{\mathsf{T}}$$

### 2.1.7  Computing $\frac{\partial J}{\partial \mathbf{W}_{ua}}$, $\frac{\partial J}{\partial \mathbf{W}_{ux}}$, and $\frac{\partial J}{\partial \mathbf{b}_u}$

From *eq.4* (see section-2), we have

$$\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} = \frac{\partial J}{\partial \Gamma_u^{\langle t \rangle}} \frac{\partial \Gamma_u^{\langle t \rangle}}{\partial \mathbf{z}_u^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \Gamma_u^{\langle t \rangle}} \circ \left(\Gamma_u^{\langle t \rangle} \circ \left(\mathbf{1}_{(n_a, 1)} - \Gamma_u^{\langle t \rangle}\right)\right)^{\mathsf{T}}$$

From *eq.3* (see section-2), using the trace-method, we have

$$\operatorname{tr}(J) = \operatorname{tr}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} \, \mathrm{d}\mathbf{z}_u^{\langle t \rangle}\right)$$

where, $\mathrm{d}\mathbf{z}_u^{\langle t \rangle}$ can be expanded as follows:

$$\mathrm{d}\mathbf{z}_u^{\langle t \rangle} = \mathrm{d}\left(\mathbf{W}_{ua}^{\mathsf{T}} \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{ux}^{\mathsf{T}} \mathbf{x}^{\langle t \rangle} + \mathbf{b}_u\right)$$

$$= \mathrm{d}\mathbf{W}_{ua}^{\mathsf{T}} \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{ua}^{\mathsf{T}} \mathrm{d}\mathbf{a}^{\langle t-1 \rangle} + \mathrm{d}\mathbf{W}_{ux}^{\mathsf{T}} \mathbf{x}^{\langle t \rangle} + \mathbf{W}_{ux}^{\mathsf{T}} \mathrm{d}\mathbf{x}^{\langle t \rangle} + \mathrm{d}\mathbf{b}_u \qquad \text{eq.13-3}$$

when differentiating w.r.t. $\mathbf{W}_{ua}$, we have $\mathrm{d}\mathbf{W}_{ux} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1 \rangle} = 0$, $\mathrm{d}\mathbf{b}_u = 0$, and $\mathrm{d}\mathbf{x}^{\langle t \rangle} = 0$. So,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} d\mathbf{W}_{ua}^{\mathsf{T}} \mathbf{a}^{\langle t-1 \rangle}\right)$$

$$= \text{tr}\left((\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}} d\mathbf{W}_{ua}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}}\right)$$

$$= \text{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}} (\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}} d\mathbf{W}_{ua}\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{W}_{ua}} = \left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}} (\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}}$$

when differentiating w.r.t. $\mathbf{W}_{ux}$, we have $d\mathbf{W}_{ua} = 0$, $d\mathbf{a}^{\langle t-1 \rangle} = 0$, $d\mathbf{b}_u = 0$, and $d\mathbf{x}^{\langle t \rangle} = 0$. So,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} d\mathbf{W}_{ux}^{\mathsf{T}} \mathbf{x}^{\langle t \rangle}\right)$$

$$= \text{tr}\left((\mathbf{x}^{\langle t \rangle})^{\mathsf{T}} d\mathbf{W}_{ux}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}}\right)$$

$$= \text{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}} (\mathbf{x}^{\langle t \rangle})^{\mathsf{T}} d\mathbf{W}_{ux}\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{W}_{ux}} = \left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}\right)^{\mathsf{T}} (\mathbf{x}^{\langle t \rangle})^{\mathsf{T}}$$

when differentiating w.r.t. $\mathbf{b}_u$, we have $d\mathbf{W}_{ua} = 0$, $d\mathbf{a}^{\langle t-1 \rangle} = 0$, $d\mathbf{W}_{ux} = 0$, and $d\mathbf{x}^{\langle t \rangle} = 0$. So,

$$dJ = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} d\mathbf{b}_u\right)$$

$$\implies \frac{\partial J}{\partial \mathbf{b}_u} = \frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}}$$

### 2.1.8 Computing $\frac{\partial J}{\partial \mathbf{W}_{fa}}$, $\frac{\partial J}{\partial \mathbf{W}_{fx}}$, and $\frac{\partial J}{\partial \mathbf{b}_f}$

From *eq.2* (see section-2), we have

$$\frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} = \frac{\partial J}{\partial \Gamma_f^{\langle t \rangle}} \frac{\partial \Gamma_f^{\langle t \rangle}}{\partial \mathbf{z}_f^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \Gamma_f^{\langle t \rangle}} \circ \left(\Gamma_f^{\langle t \rangle} \circ (\mathbf{1}_{(n_a, 1)} - \Gamma_f^{\langle t \rangle})\right)^{\mathsf{T}}$$

From *eq.1* (see section-2), using the trace-method, we have

$$\text{tr}(J) = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} d\mathbf{z}_f^{\langle t \rangle}\right)$$

where, $d\mathbf{z}_f^{\langle t \rangle}$ can be expanded as follows:

$$d\mathbf{z}_f^{\langle t \rangle} = d\left(\mathbf{W}_{fa}^{\mathsf{T}} \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{fx}^{\mathsf{T}} \mathbf{x}^{\langle t \rangle} + \mathbf{b}_f\right)$$

$$= d\mathbf{W}_{fa}^{\mathsf{T}} \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{fa}^{\mathsf{T}} d\mathbf{a}^{\langle t-1 \rangle} + d\mathbf{W}_{fx}^{\mathsf{T}} \mathbf{x}^{\langle t \rangle} + \mathbf{W}_{fx}^{\mathsf{T}} d\mathbf{x}^{\langle t \rangle} + d\mathbf{b}_f \qquad \text{eq.13-4}$$

when differentiating w.r.t. $\mathbf{W}_{fa}$, we have $\mathrm{d}\mathbf{W}_{fx} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1 \rangle} = 0$, $\mathrm{d}\mathbf{b}_f = 0$, and $\mathrm{d}\mathbf{x}^{\langle t \rangle} = 0$. So,

$$
\begin{aligned}
\mathrm{d}J &= \mathrm{tr}\left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \, \mathrm{d}\mathbf{W}_{fa}^{\mathsf{T}} \, \mathbf{a}^{\langle t-1 \rangle} \right) \\
&= \mathrm{tr}\left( (\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}} \, \mathrm{d}\mathbf{W}_{fa} \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} \right) \\
&= \mathrm{tr}\left( \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} (\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}} \, \mathrm{d}\mathbf{W}_{fa} \right) \\
\implies \frac{\partial J}{\partial \mathbf{W}_{fa}} &= \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} (\mathbf{a}^{\langle t-1 \rangle})^{\mathsf{T}}
\end{aligned}
$$

when differentiating w.r.t. $\mathbf{W}_{fx}$, we have $\mathrm{d}\mathbf{W}_{fa} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1 \rangle} = 0$, $\mathrm{d}\mathbf{b}_f = 0$, and $\mathrm{d}\mathbf{x}^{\langle t \rangle} = 0$. So,

$$
\begin{aligned}
\mathrm{d}J &= \mathrm{tr}\left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \, \mathrm{d}\mathbf{W}_{fx}^{\mathsf{T}} \, \mathbf{x}^{\langle t \rangle} \right) \\
&= \mathrm{tr}\left( (\mathbf{x}^{\langle t \rangle})^{\mathsf{T}} \, \mathrm{d}\mathbf{W}_{fx} \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} \right) \\
&= \mathrm{tr}\left( \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} (\mathbf{x}^{\langle t \rangle})^{\mathsf{T}} \, \mathrm{d}\mathbf{W}_{fx} \right) \\
\implies \frac{\partial J}{\partial \mathbf{W}_{fx}} &= \left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \right)^{\mathsf{T}} (\mathbf{x}^{\langle t \rangle})^{\mathsf{T}}
\end{aligned}
$$

when differentiating w.r.t. $\mathbf{b}_f$, we have $\mathrm{d}\mathbf{W}_{fa} = 0$, $\mathrm{d}\mathbf{a}^{\langle t-1 \rangle} = 0$, $\mathrm{d}\mathbf{W}_{fx} = 0$, and $\mathrm{d}\mathbf{x}^{\langle t \rangle} = 0$. So,

$$
\begin{aligned}
\mathrm{d}J &= \mathrm{tr}\left( \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \, \mathrm{d}\mathbf{b}_f \right) \\
\implies \frac{\partial J}{\partial \mathbf{b}_f} &= \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}}
\end{aligned}
$$

### 2.1.9 Computing $\frac{\partial J}{\partial \mathbf{W}_{ca}}$, $\frac{\partial J}{\partial \mathbf{W}_{cx}}$, and $\frac{\partial J}{\partial \mathbf{b}_c}$

From *eq.8* (see section-2), we have

$$
\begin{aligned}
\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{\langle t \rangle}} \frac{\partial \tilde{\mathbf{c}}^{\langle t \rangle}}{\partial \mathbf{z}_c^{\langle t \rangle}} \\
&= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{\langle t \rangle}} \circ \left[ \mathbf{1}_{(n_a,1)} - \tanh^2(\mathbf{z}_c^{\langle t \rangle}) \right]^{\mathsf{T}}
\end{aligned}
$$

From *eq.7* (see section-2), using the trace-method, we have

$$\text{tr}(J) = \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \, d\mathbf{z}_c^{\langle t \rangle}\right)$$

where, $d\mathbf{z}_c^{\langle t \rangle}$ can be expanded as follows:

$$\begin{aligned}
d\mathbf{z}_c^{\langle t \rangle} &= d\left(\mathbf{W}_{ca}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{cx}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{b}_c\right) \\
&= d\mathbf{W}_{ca}^\intercal \, \mathbf{a}^{\langle t-1 \rangle} + \mathbf{W}_{ca}^\intercal \, d\mathbf{a}^{\langle t-1 \rangle} + d\mathbf{W}_{cx}^\intercal \, \mathbf{x}^{\langle t \rangle} + \mathbf{W}_{cx}^\intercal \, d\mathbf{x}^{\langle t \rangle} + d\mathbf{b}_c \qquad \text{eq.13-5}
\end{aligned}$$

when differentiating w.r.t. $\mathbf{W}_{ca}$, we have $d\mathbf{W}_{cx} = 0$, $d\mathbf{a}^{\langle t-1 \rangle} = 0$, $d\mathbf{b}_c = 0$, and $d\mathbf{x}^{\langle t \rangle} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \, d\mathbf{W}_{ca}^\intercal \, \mathbf{a}^{\langle t-1 \rangle}\right) \\
&= \text{tr}\left(\left(\mathbf{a}^{\langle t-1 \rangle}\right)^\intercal \, d\mathbf{W}_{ca} \left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal\right) \\
&= \text{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal \left(\mathbf{a}^{\langle t-1 \rangle}\right)^\intercal \, d\mathbf{W}_{ca}\right) \\
\implies \frac{\partial J}{\partial \mathbf{W}_{ca}} &= \left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal \left(\mathbf{a}^{\langle t-1 \rangle}\right)^\intercal
\end{aligned}$$

when differentiating w.r.t. $\mathbf{W}_{cx}$, we have $d\mathbf{W}_{ca} = 0$, $d\mathbf{a}^{\langle t-1 \rangle} = 0$, $d\mathbf{b}_c = 0$, and $d\mathbf{x}^{\langle t \rangle} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \, d\mathbf{W}_{cx}^\intercal \, \mathbf{x}^{\langle t \rangle}\right) \\
&= \text{tr}\left(\left(\mathbf{x}^{\langle t \rangle}\right)^\intercal \, d\mathbf{W}_{cx} \left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal\right) \\
&= \text{tr}\left(\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal \left(\mathbf{x}^{\langle t \rangle}\right)^\intercal \, d\mathbf{W}_{cx}\right) \\
\implies \frac{\partial J}{\partial \mathbf{W}_{cx}} &= \left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}\right)^\intercal \left(\mathbf{x}^{\langle t \rangle}\right)^\intercal
\end{aligned}$$

when differentiating w.r.t. $\mathbf{b}_c$, we have $d\mathbf{W}_{ca} = 0$, $d\mathbf{a}^{\langle t-1 \rangle} = 0$, $d\mathbf{W}_{cx} = 0$, and $d\mathbf{x}^{\langle t \rangle} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr}\left(\frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \, d\mathbf{b}_c\right) \\
\implies \frac{\partial J}{\partial \mathbf{b}_c} &= \frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}}
\end{aligned}$$

### 2.1.10 Computing $\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}}$, and $\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}}$

Each of these derivatives has four components, described as follows

**Comp-1** flows-in as derivative from $\Gamma_o^{\langle t \rangle}$

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} \frac{\partial \mathbf{z}_o^{\langle t \rangle}}{\partial \mathbf{a}^{\langle t-1 \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} \mathbf{W}_{oa}^{\intercal}$$

$$\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} \frac{\partial \mathbf{z}_o^{\langle t \rangle}}{\partial \mathbf{x}^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_o^{\langle t \rangle}} \mathbf{W}_{ox}^{\intercal}$$

**Comp-2** flows in as derivative from $\tilde{\mathbf{c}}^{\langle t \rangle}$

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \frac{\partial \mathbf{z}_c^{\langle t \rangle}}{\partial \mathbf{a}^{\langle t-1 \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \mathbf{W}_{ca}^{\intercal}$$

$$\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \frac{\partial \mathbf{z}_c^{\langle t \rangle}}{\partial \mathbf{x}^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_c^{\langle t \rangle}} \mathbf{W}_{cx}^{\intercal}$$

**Comp-3** flows-in as derivative from $\Gamma_u^{\langle t \rangle}$

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} \frac{\partial \mathbf{z}_u^{\langle t \rangle}}{\partial \mathbf{a}^{\langle t-1 \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} \mathbf{W}_{ua}^{\intercal}$$

$$\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} \frac{\partial \mathbf{z}_u^{\langle t \rangle}}{\partial \mathbf{x}^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_u^{\langle t \rangle}} \mathbf{W}_{ux}^{\intercal}$$

**Comp-4** flows-in as derivative from $\Gamma_f^{\langle t \rangle}$

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \frac{\partial \mathbf{z}_f^{\langle t \rangle}}{\partial \mathbf{a}^{\langle t-1 \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \mathbf{W}_{fa}^{\intercal}$$

$$\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \frac{\partial \mathbf{z}_f^{\langle t \rangle}}{\partial \mathbf{x}^{\langle t \rangle}}$$

$$= \frac{\partial J}{\partial \mathbf{z}_f^{\langle t \rangle}} \mathbf{W}_{fx}^{\intercal}$$

These four components are then added to compute the true derivative, i.e.

$$\frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}} = \frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}}\bigg|_{\text{Comp-1}} + \frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}}\bigg|_{\text{Comp-2}} + \frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}}\bigg|_{\text{Comp-3}} + \frac{\partial J}{\partial \mathbf{a}^{\langle t-1 \rangle}}\bigg|_{\text{Comp-4}}$$

(see $\oplus$ labeled as *3*, *4*, and *5* in figure-2)

$$\frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}} = \frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}}\bigg|_{\text{Comp-1}} + \frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}}\bigg|_{\text{Comp-2}} + \frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}}\bigg|_{\text{Comp-3}} + \frac{\partial J}{\partial \mathbf{x}^{\langle t \rangle}}\bigg|_{\text{Comp-4}}$$

(see $\oplus$ labeled as *6*, *7*, and *8* in figure-2)

# 3  Gradient or Jacobian?

In the above derivations, we have used the numerator layout while performing matrix-derivatives. One of the consequences of this decision is that the derivatives that we have computed are in-fact jacobians and not gradients. Fortunately, gradients are just transpose of jacobians.

# 4  Appendix-A:

In this section we will derive an expression for the derivative of a Hadamard product of two vectors. Let $\mathbf{x}$ and $\mathbf{y}$ be two vectors, and let $\mathbf{w} = \mathbf{x} \circ \mathbf{y}$ such that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} ; \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} ; \qquad \mathbf{w} = \mathbf{x} \circ \mathbf{y} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \\ \vdots \\ x_n y_n \end{bmatrix}$$

then, we have

$$
\begin{aligned}
\frac{\partial \mathbf{w}}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} & \cdots & \frac{\partial}{\partial x_n} \end{bmatrix} \otimes \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \\ \vdots \\ x_n y_n \end{bmatrix} \\
&= \begin{bmatrix} \frac{\partial(x_1 y_1)}{\partial x_1} & \frac{\partial(x_1 y_1)}{\partial x_2} & \cdots & \frac{\partial(x_1 y_1)}{\partial x_n} \\ \frac{\partial(x_2 y_2)}{\partial x_1} & \frac{\partial(x_2 y_2)}{\partial x_2} & \cdots & \frac{\partial(x_2 y_2)}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial(x_n y_n)}{\partial x_1} & \frac{\partial(x_n y_n)}{\partial x_2} & \cdots & \frac{\partial(x_n y_n)}{\partial x_n} \end{bmatrix} \\
&= \begin{bmatrix} y_1 & 0 & \cdots & 0 \\ 0 & y_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & y_n \end{bmatrix} \\
&= \text{diag}\,(\mathbf{y})
\end{aligned}
\tag{1}
$$