$$S = QK^T$$

$$P = \text{softmax}(S)$$

$$O = PV$$

$$\begin{bmatrix} [P_{11} \ P_{12} \ P_{13}] \end{bmatrix} \begin{bmatrix} V_{11} & V_{12} & V_{13} & V_{14} \\ V_{21} & V_{22} & V_{23} & V_{24} \\ V_{31} & V_{32} & V_{33} & V_{34} \end{bmatrix}$$

1st row
$$= \begin{bmatrix} (P_{11}V_{11} + P_{12}V_{21} + P_{13}V_{31}) & \begin{pmatrix} P_{11}V_{12} + P_{12}V_{22} \\ + P_{13}V_{32} \end{pmatrix} & \begin{pmatrix} P_{11}V_{13} + P_{12}V_{23} \\ + P_{13}V_{23} \end{pmatrix} & \begin{pmatrix} P_{11}V_{14} + P_{12}V_{24} \\ + P_{13}V_{34} \end{pmatrix} \end{bmatrix}$$

$$O_1 = P_{11}\begin{bmatrix} V_{1:} \end{bmatrix} + P_{12}\begin{bmatrix} V_{2:} \end{bmatrix} + P_{13}\begin{bmatrix} V_{3:} \end{bmatrix}$$

$$\boxed{O_i = \sum_j P_{ij}(V_j)}$$

$$Y = XW$$

$$\frac{d\phi}{dX} = dX = \frac{d\phi}{dY}W^T$$

$$\frac{d\phi}{dW} = dW = X^T\frac{d\phi}{dY}$$

When $O = PV$

$$dP = dO \cdot V^T$$

$$dV = P^T dO$$

$$dV_i = \sum_j (P^T)_{ij} \cdot dO_j$$

$$dV_i = \sum_j P_{ji} \, dO_j$$

$$dP_i = \sum_j dO_{ij} (V^T)_j$$

___

$P_{i:} = \text{softmax}(S_{i:})$

For a vector $x, y$, s.t. $y = \text{softmax}(x)$

Jacobian $= \boxed{\text{diag}(y) - y \cdot y^T} \longrightarrow \underline{\text{Symmetric}}$

$dx = dS_{i:} = \left( \text{diag}(P_{i:}) - P_{i:} P_{i:}^T \right)$

$\dfrac{d\phi}{dS} = \dfrac{d\phi}{dP} \cdot \dfrac{dP}{dS} \longrightarrow$ Row convention

$\Rightarrow$ This is in column convention

$\Rightarrow \left( \dfrac{d\phi}{dS} \right)^T = \dfrac{dP}{dS} \cdot \left( \dfrac{d\phi}{dP} \right)^T$

$\underset{\text{Symmetric}}{\downarrow} \qquad = dP_{i:}$

$$dS_{i:} = \left( \text{diag}(P_{i:}) - P_{i:} P_{i:}^T \right) dP_{i:} \;\Leftarrow$$

$$= \left( P_{i:} \odot dP_{i:} \right) - \underbrace{\left( P_{i:}^T \, dP_{i:} \right)}_{\text{Strange but equivalent}} P_{i:}$$

<u>Define</u>

$$D_i = P_{i:}^T \, dP_{i:}$$

$$D_i = \sum_j \left( P_{i:}^T \right)_{ij} \left( dP_{i:} \right)_j$$

$$= \sum_j \frac{e^{q_i^T k_j}}{L_i} \cdot \underbrace{dO_i^T \cdot V_j}_{|\times|} = dO_i^T \sum_j \frac{e^{q_i^T k_j}}{L_i} V_j$$

$$= dO_i^T \cdot O_i$$

$$\Rightarrow \boxed{D_i = dO_i^T \, O_i}$$

$$\boxed{dS_{i:} = P_{i:} \odot dP_{i:} - D_i \, P_{i:}}$$

$$dS_{ij} = \left( P_{ij} \times dP_{ij} \right) - D_i P_{ij}$$

$$\Rightarrow \boxed{dS_{ij} = P_{ij}\left(dP_{ij} - D_i\right)}$$

Now, $\quad S_{ij} = q_i^T k_j \Rightarrow S = Q K^T$

$$\boxed{\begin{array}{l} dQ = dS . K \\ dK = dS^T . Q \end{array}}$$

$$\begin{array}{l}
dq_i = dS_{i:} K = \sum_j dS_{ij} k_j \\[2mm]
\qquad = \sum_j P_{ij}\left(dP_{ij} - D_i\right) k_j \\[2mm]
\qquad = \sum_j \dfrac{e^{q_i^T k_j}}{L_i}\left(do_i^T v_j - D_i\right) k_j
\end{array}$$

Similarly

$$dk_j = \underset{\underset{\text{Transpose}}{\text{After}}}{\sum_i dS_{ij}} q_i = \sum_i P_{ij}\left(dP_{ij} - D_i\right) q_i$$
$$\qquad\qquad\qquad = \sum_i \dfrac{e^{q_i^T k_j}}{L_i}\left(do_i^T v_j - D_i\right) q_i$$