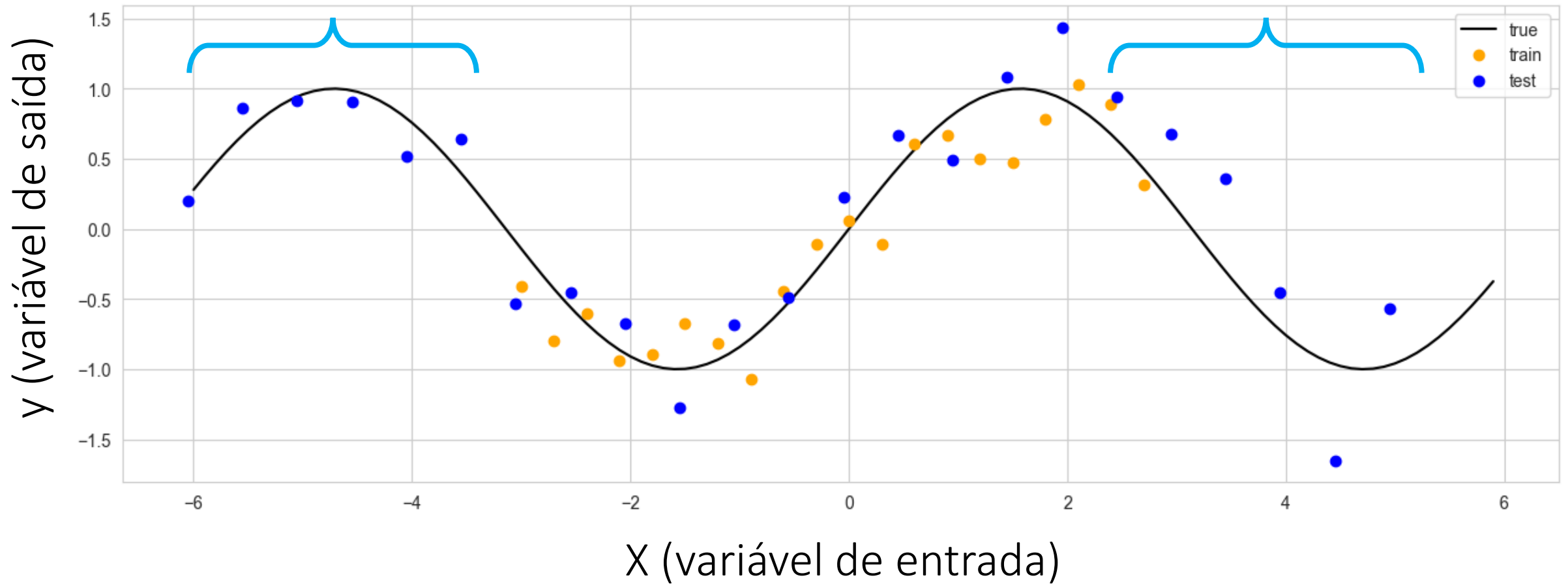
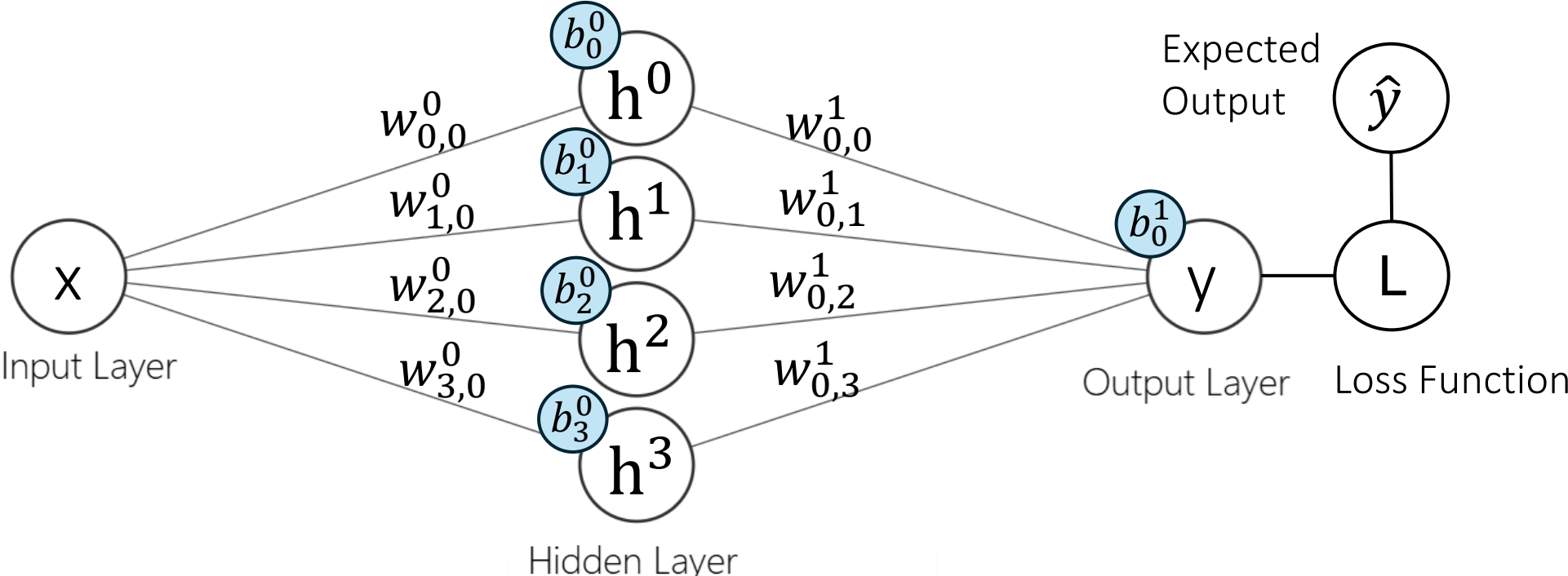


Derivação do *Backpropagation* para
uma MLP com uma camada oculta
na tarefa de regressão univariada

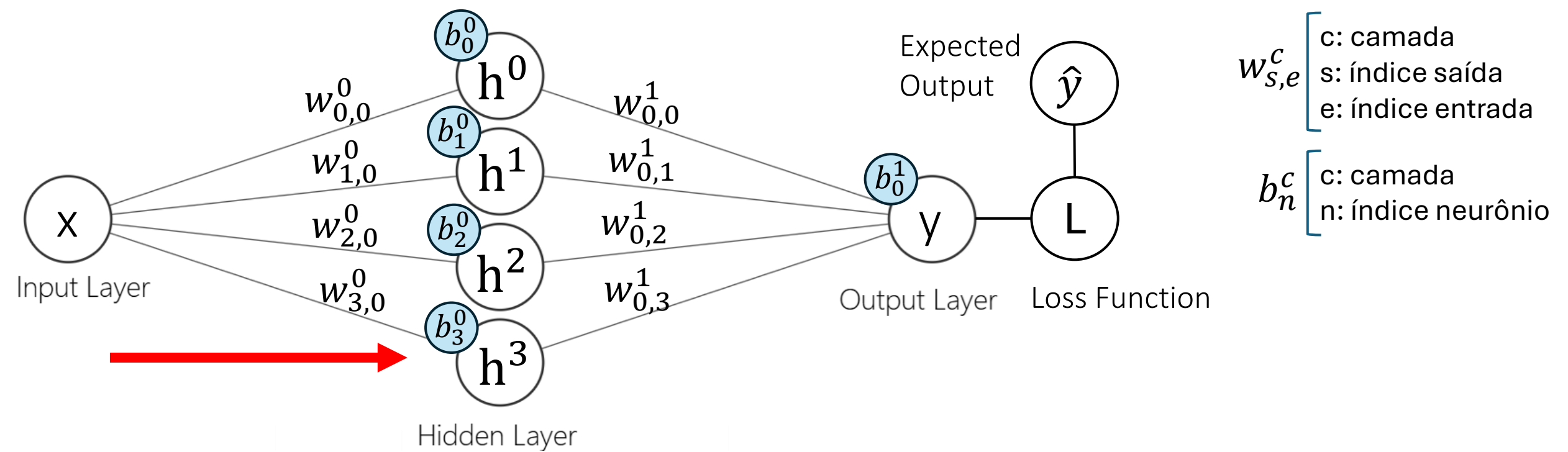
Tarefa



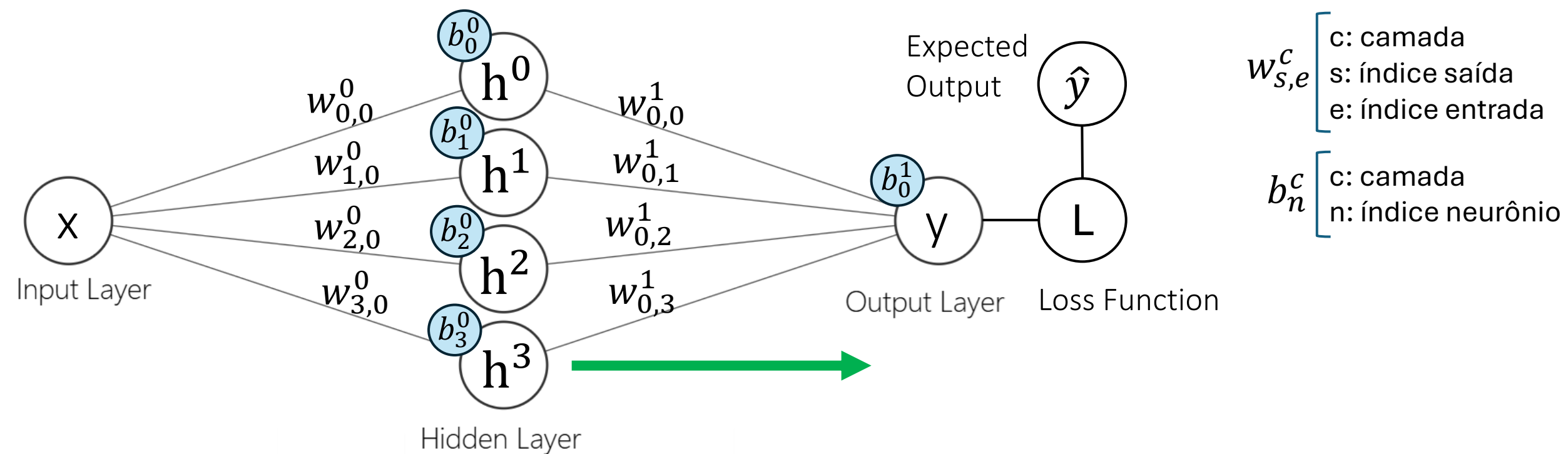


$w_{s,e}^c$ $\left\{ \begin{array}{l} c: \text{camada} \\ s: \text{índice saída} \\ e: \text{índice entrada} \end{array} \right.$

b_n^c $\left\{ \begin{array}{l} c: \text{camada} \\ n: \text{índice neurônio} \end{array} \right.$

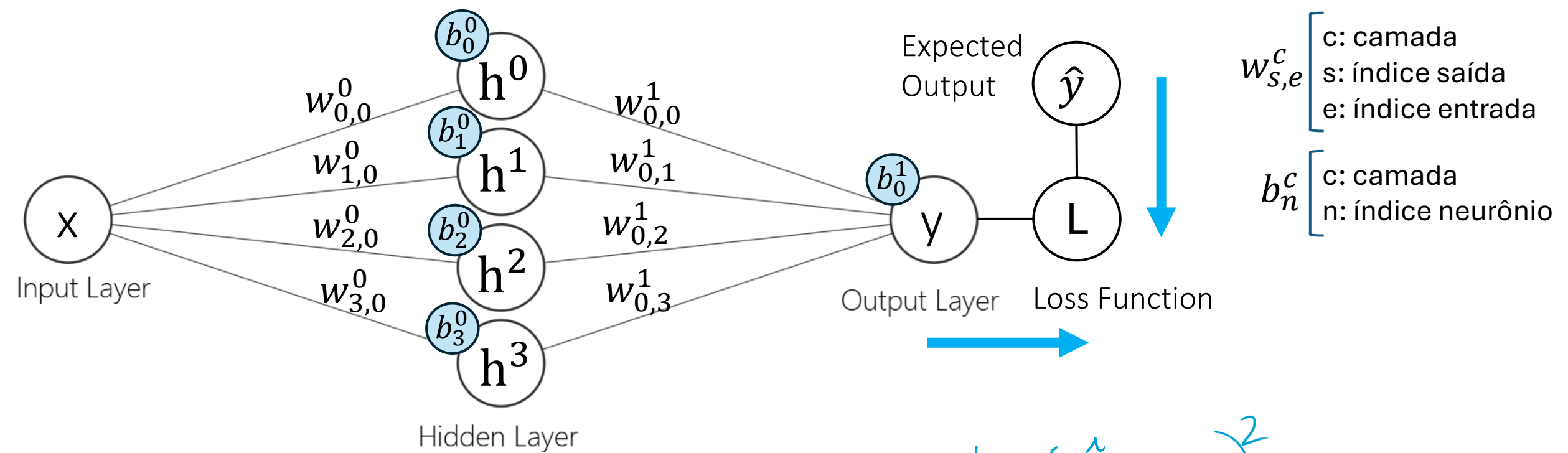


$$\sigma \left(\begin{bmatrix} w_{0,0}^0 \\ w_{1,0}^0 \\ w_{2,0}^0 \\ w_{3,0}^0 \end{bmatrix} [x_i] + \begin{bmatrix} b_0^0 \\ b_1^0 \\ b_2^0 \\ b_3^0 \end{bmatrix} \right) = \begin{bmatrix} h^0 \\ h^1 \\ h^2 \\ h^3 \end{bmatrix}$$



como a tarefa é de regressão, a função de ativação da última camada será a identidade

$$\sigma \left(\begin{bmatrix} w_{0,0}^1 & w_{0,1}^1 & w_{0,2}^1 & w_{0,3}^1 \end{bmatrix} \begin{bmatrix} h^0 \\ h^1 \\ h^2 \\ h^3 \end{bmatrix} + [b_0^1] \right) = [y_i]$$

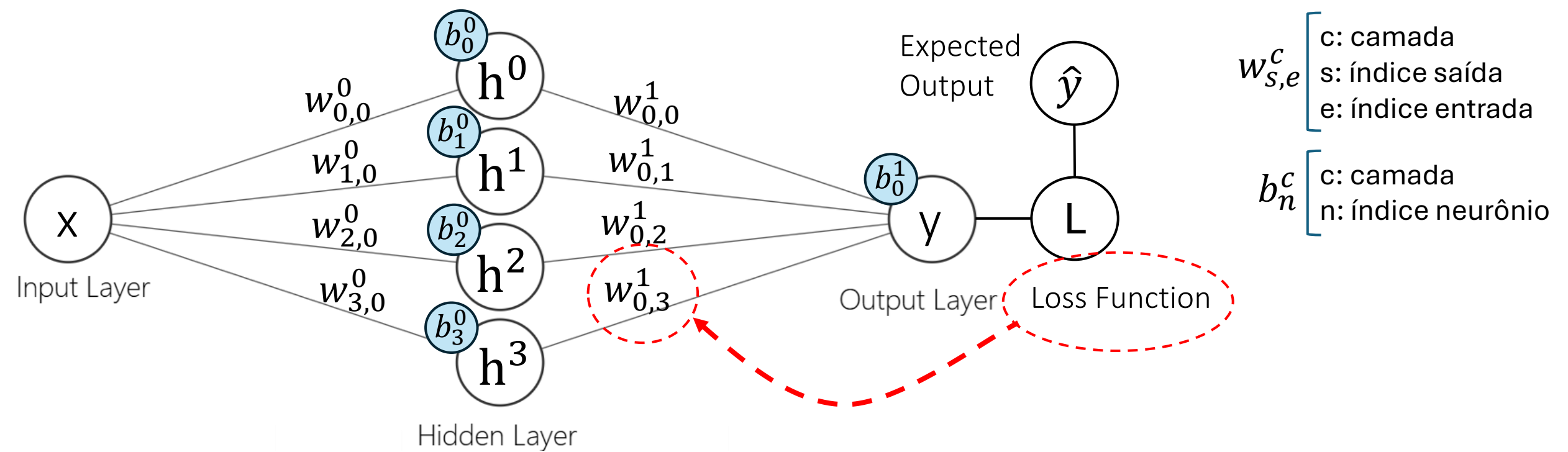


para amostra i :

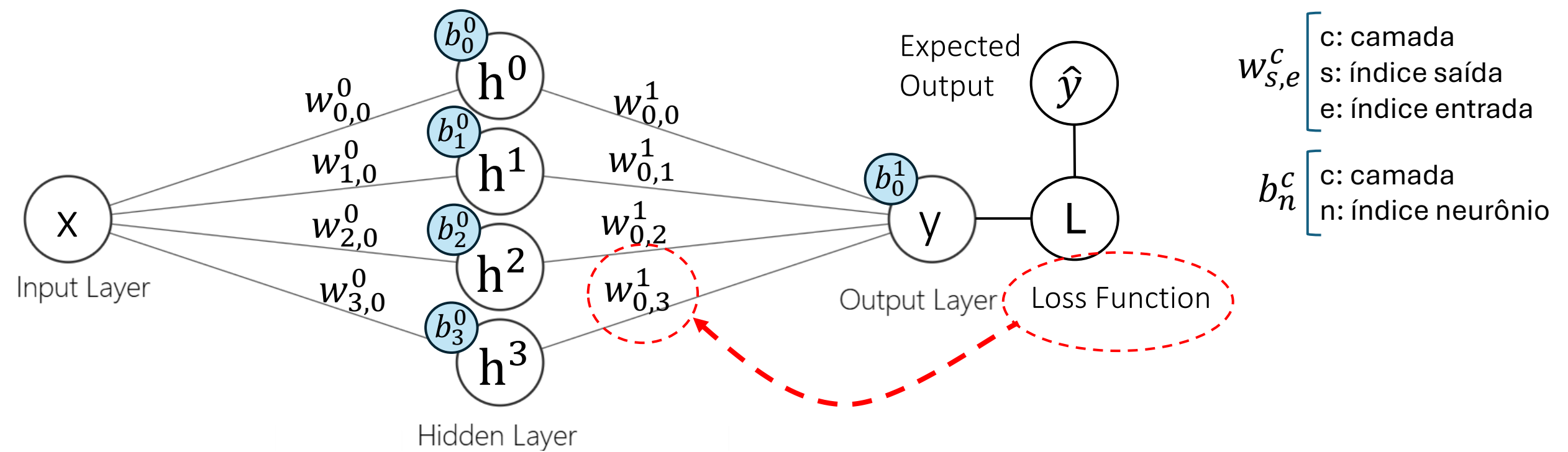
$$L_i = (\hat{y}_i - y_i)^2$$

MSE para dataset:

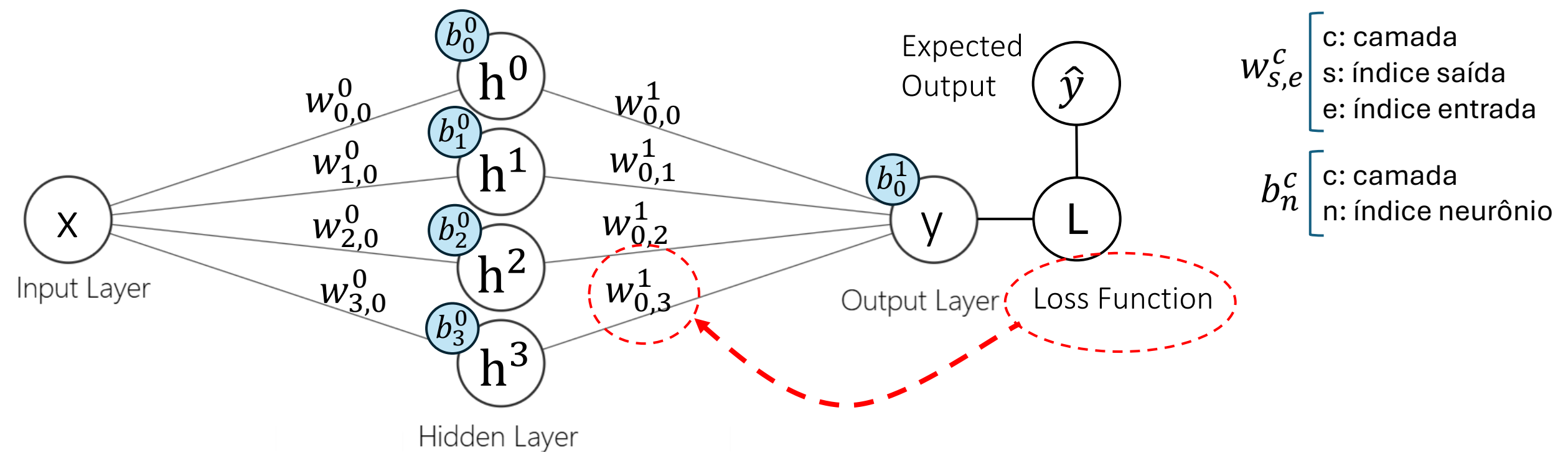
$$L = \frac{1}{N} \sum_{i=0}^N (\hat{y}_i - y_i)^2$$



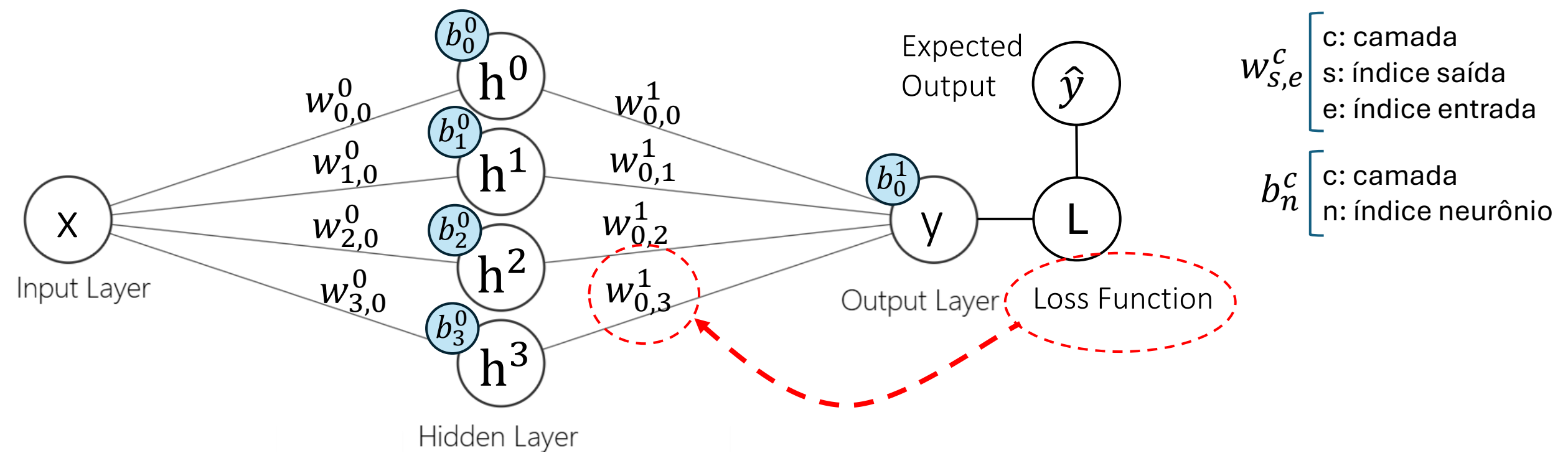
$$\frac{\partial L}{\partial w_{0,3}^1}$$



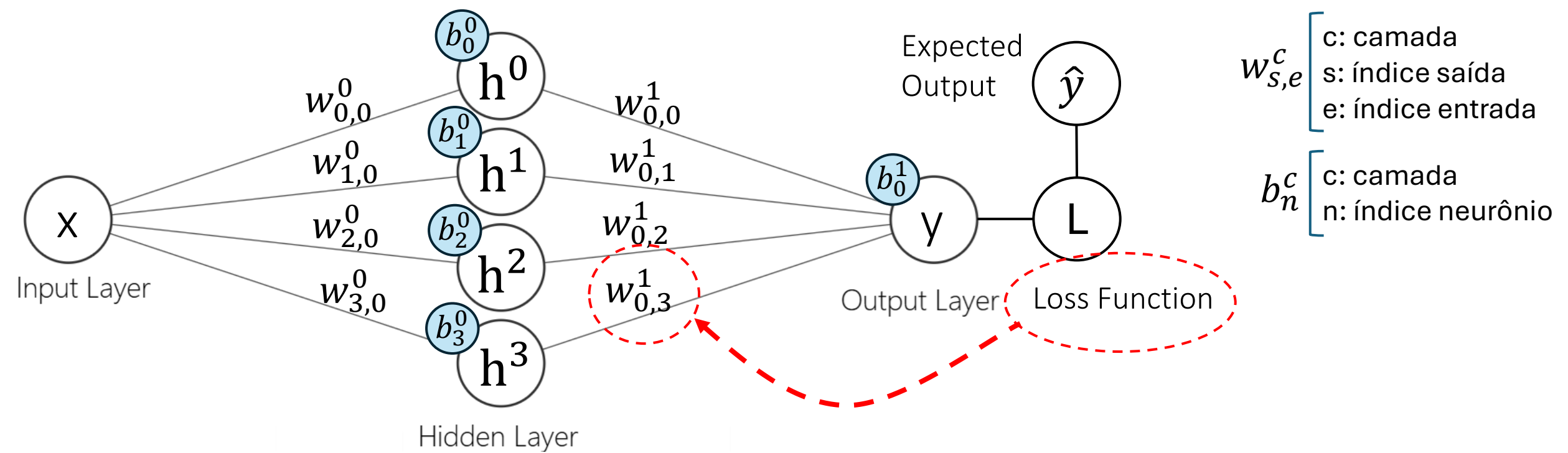
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2$$



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2$$



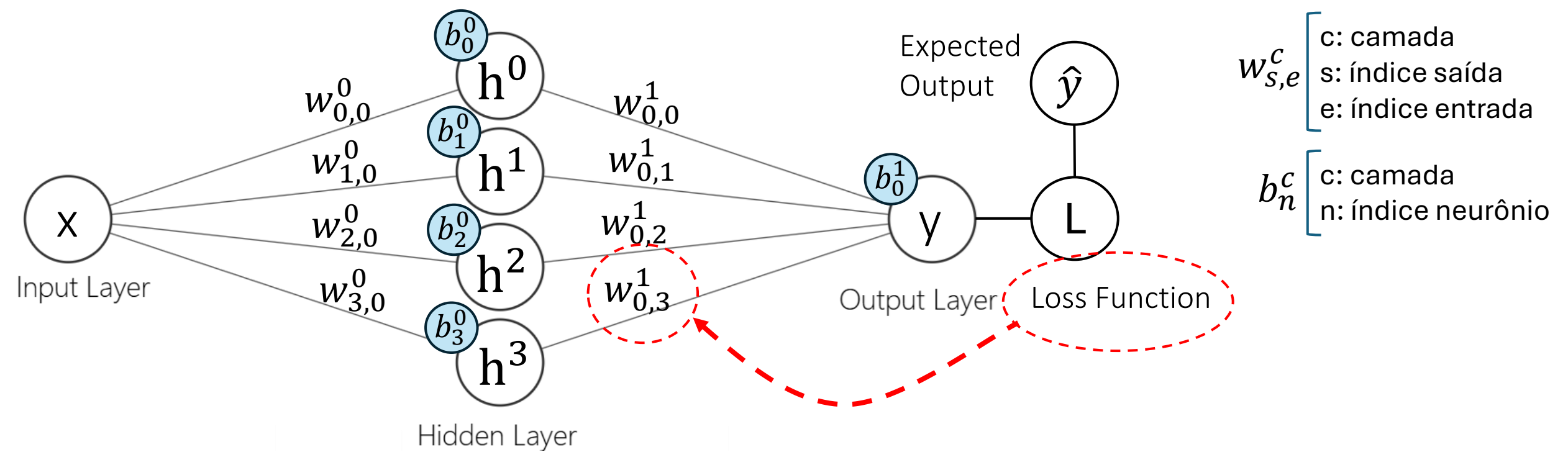
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{0,3}^1} y_i$$



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{0,3}^1} y_i$$

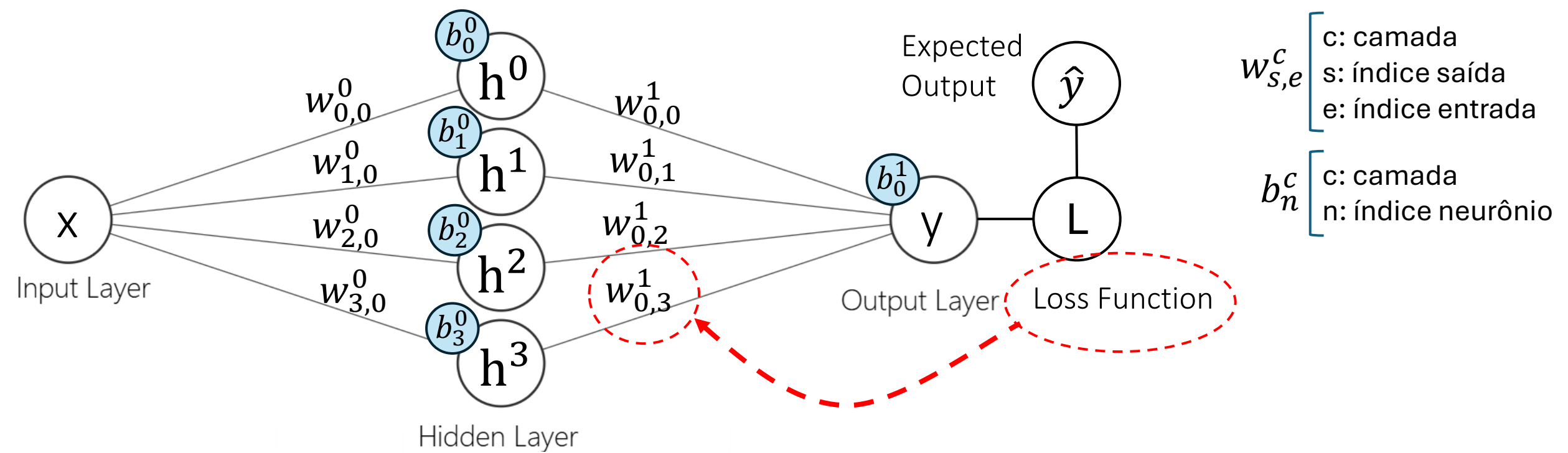
$$\frac{\partial}{\partial w_{0,3}^1} y_i = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1]$$

lembrando que a função de ativação é identidade na última camada



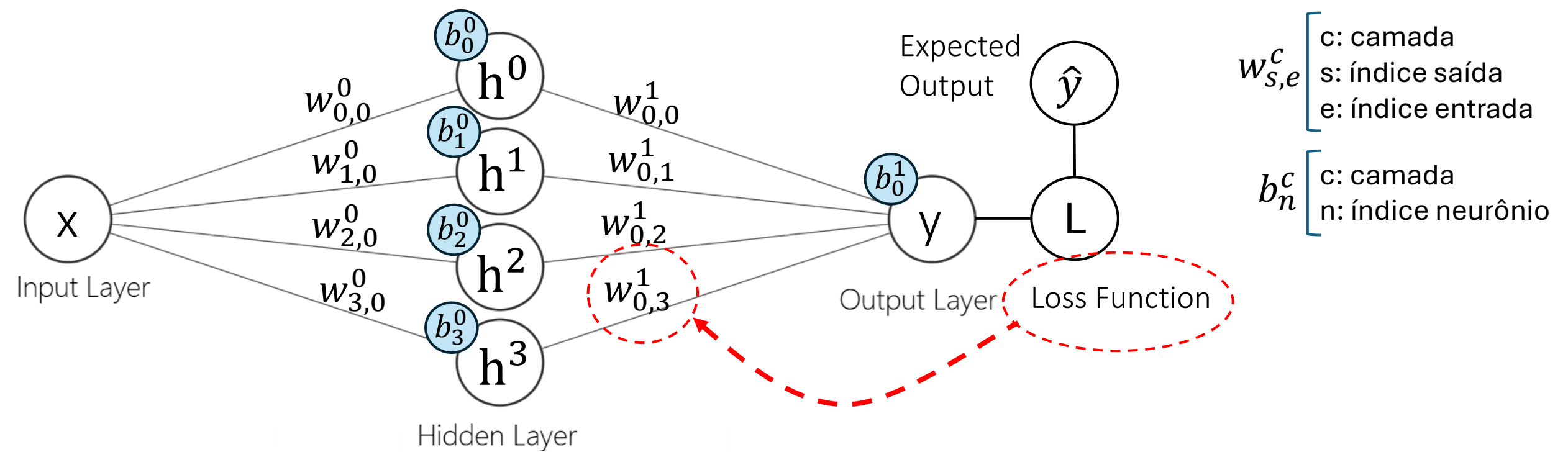
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{0,3}^1} y_i$$

$$\frac{\partial}{\partial w_{0,3}^1} y_i = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1]$$

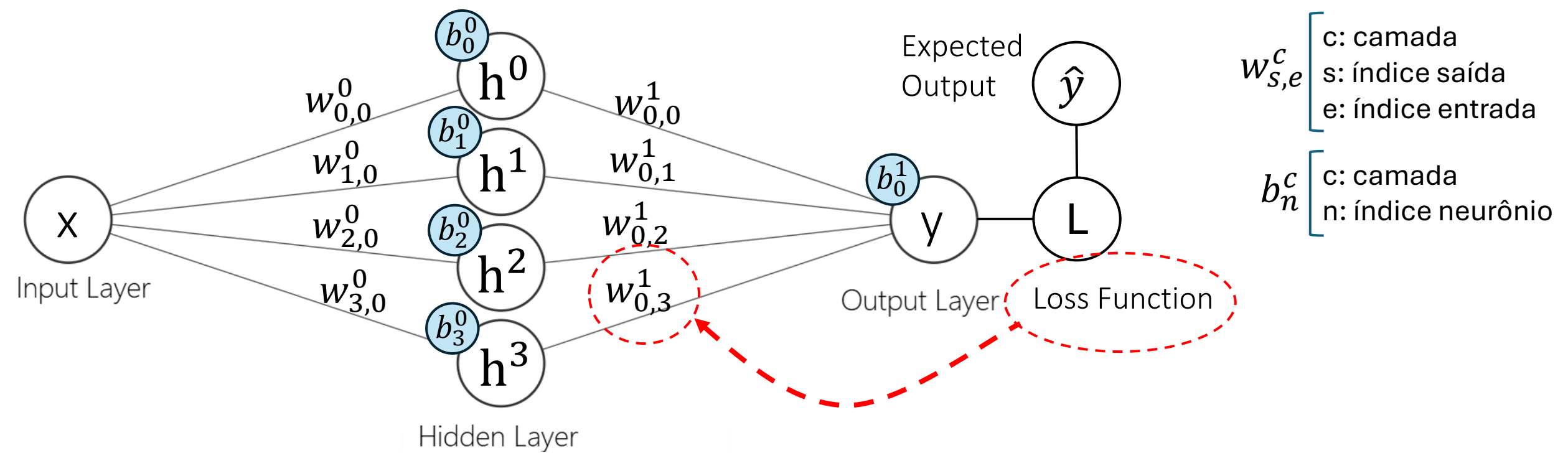


$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{0,3}^1} y_i$$

$$\frac{\partial}{\partial w_{0,3}^1} y_i = h^3$$

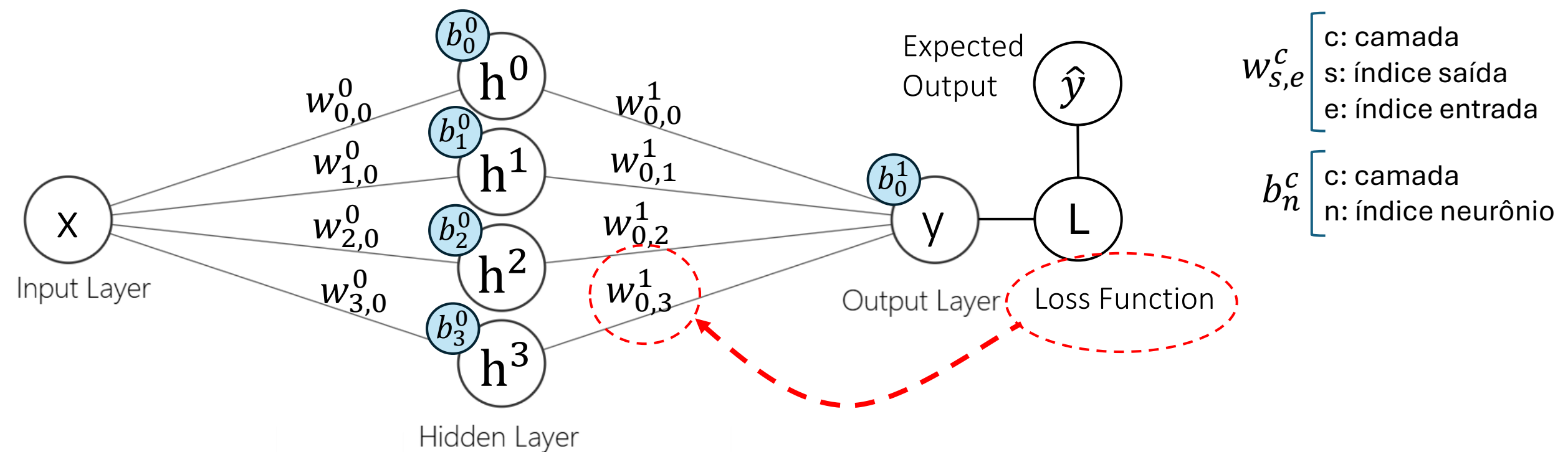


$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial w_{0,3}^1} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$

Tire um momento para refletir sobre a intuição da equação

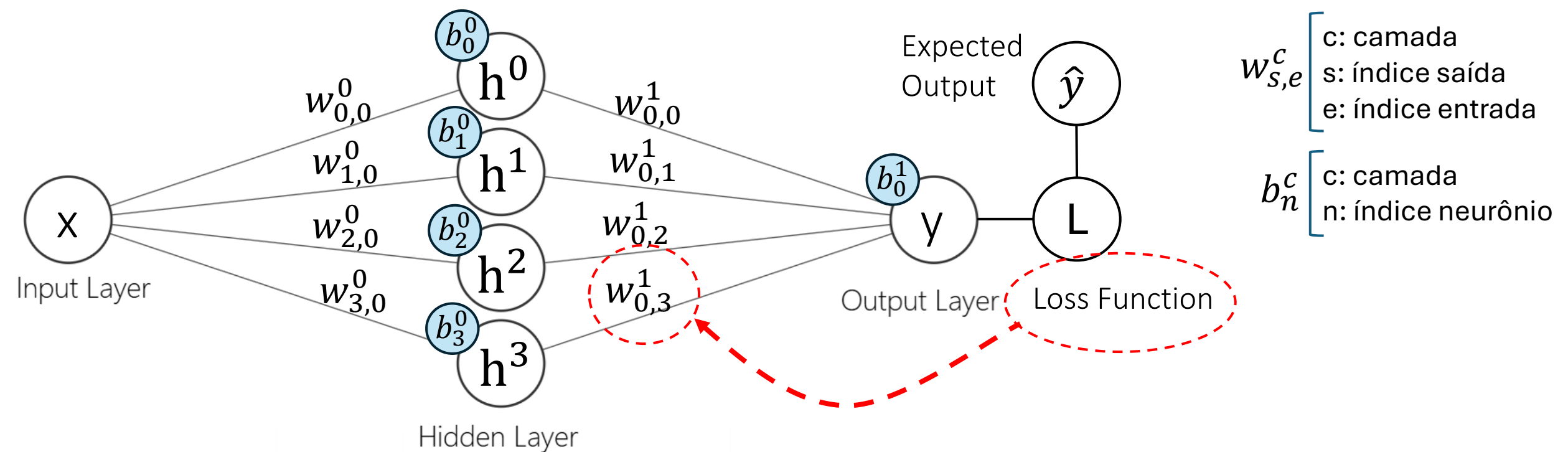


$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$

A blue arrow points from the h^3 term in the equation to the h^3 node in the diagram above.

- $w_{0,3}^1$ influencia L via y_i

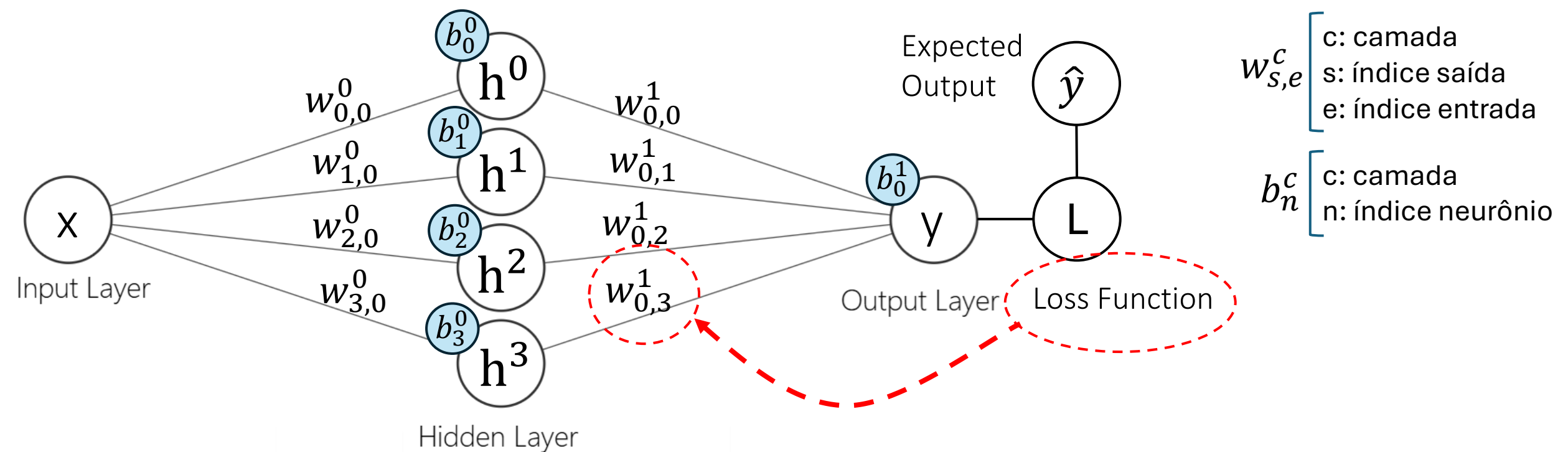
Tire um momento para refletir sobre a intuição da equação



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N \underbrace{2(y_i - \hat{y}_i)}_{\text{Loss}} h^3$$

- $w_{0,3}^1$ influencia L via y_i
- mas o impacto na *loss* também depende de \hat{y}_i

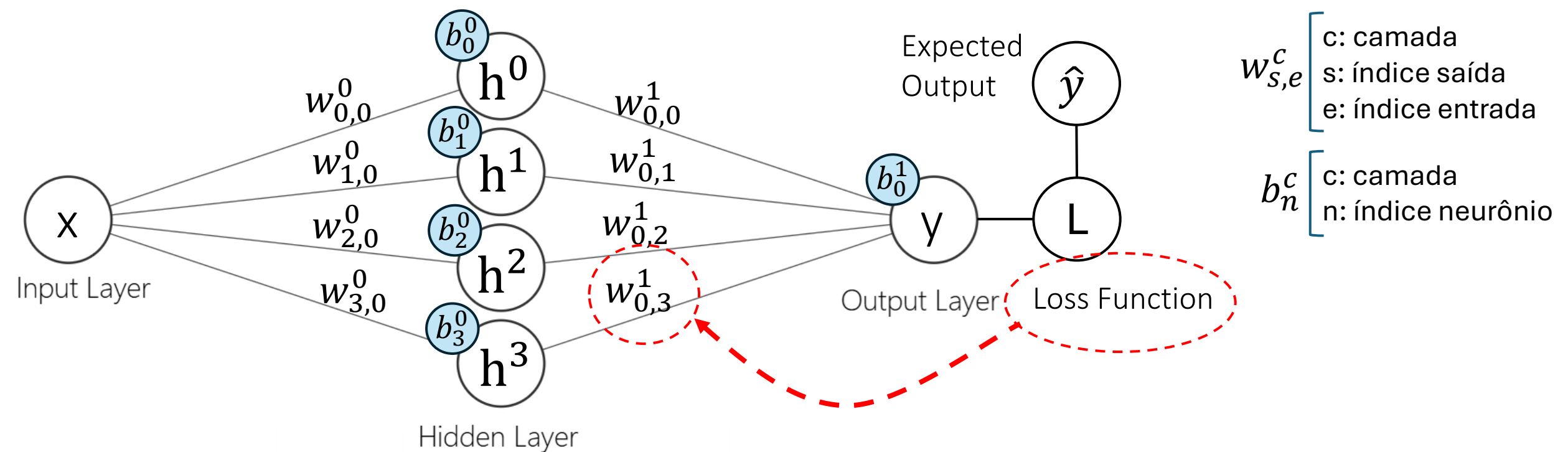
Tire um momento para refletir sobre a intuição da equação



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$

- $w_{0,3}^1$ influencia L via y_i
- mas o impacto na *loss* também depende de \hat{y}_i
- h^3 define o quanto $w_{0,3}^1$ vai mudar y_i

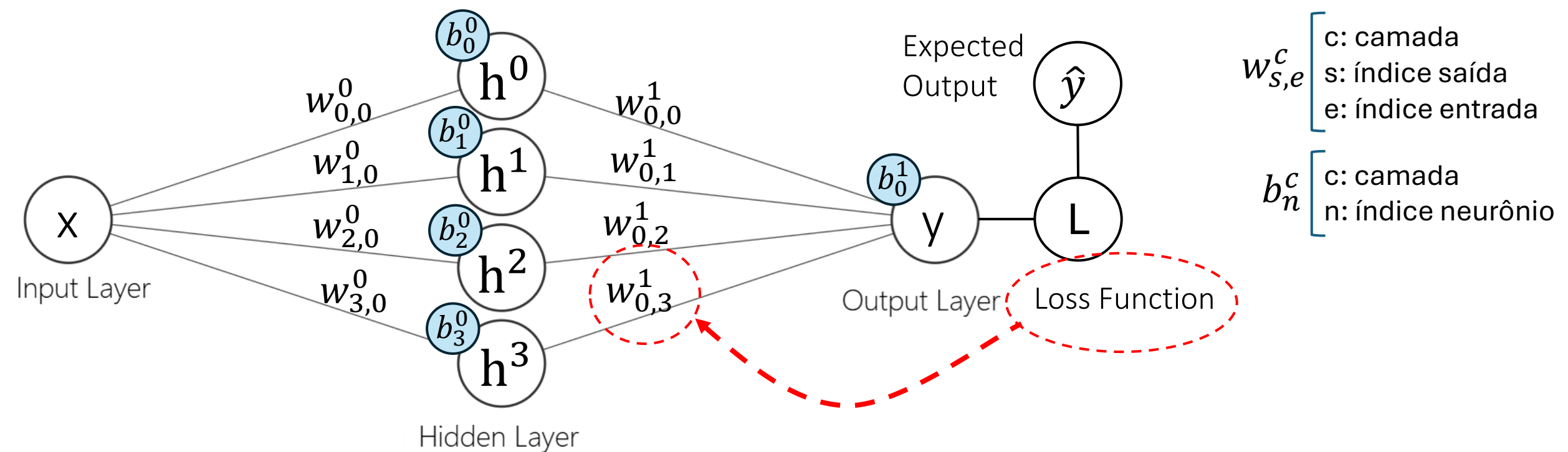
Tire um momento para refletir sobre a intuição da equação



$$\frac{\partial L}{\partial w^1_{0,3}} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$

Tire um momento para refletir sobre a intuição da equação

- $w^1_{0,3}$ influencia L via y_i
- mas o impacto na *loss* também depende de \hat{y}_i
- h^3 define o quanto $w^1_{0,3}$ vai mudar y_i
- a modificação em $w^1_{0,3}$ será a média das modificações por amostra de treino.



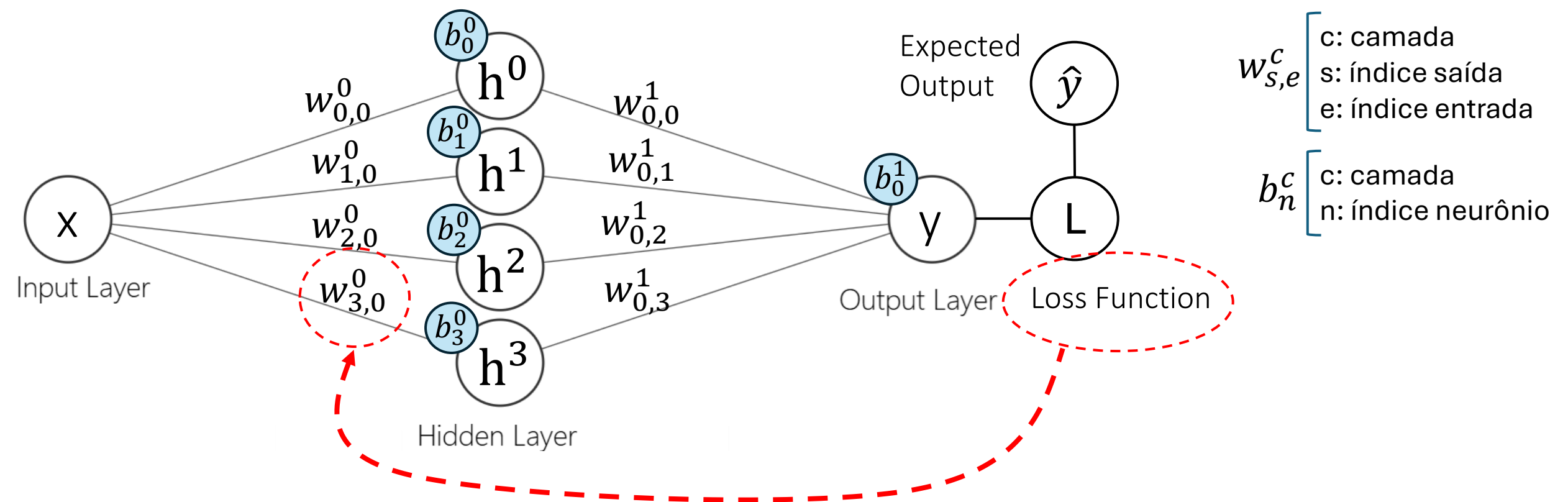
$$\frac{\partial L}{\partial w_{0,k}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^k$$

$$\frac{\partial L}{\partial b_0^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$

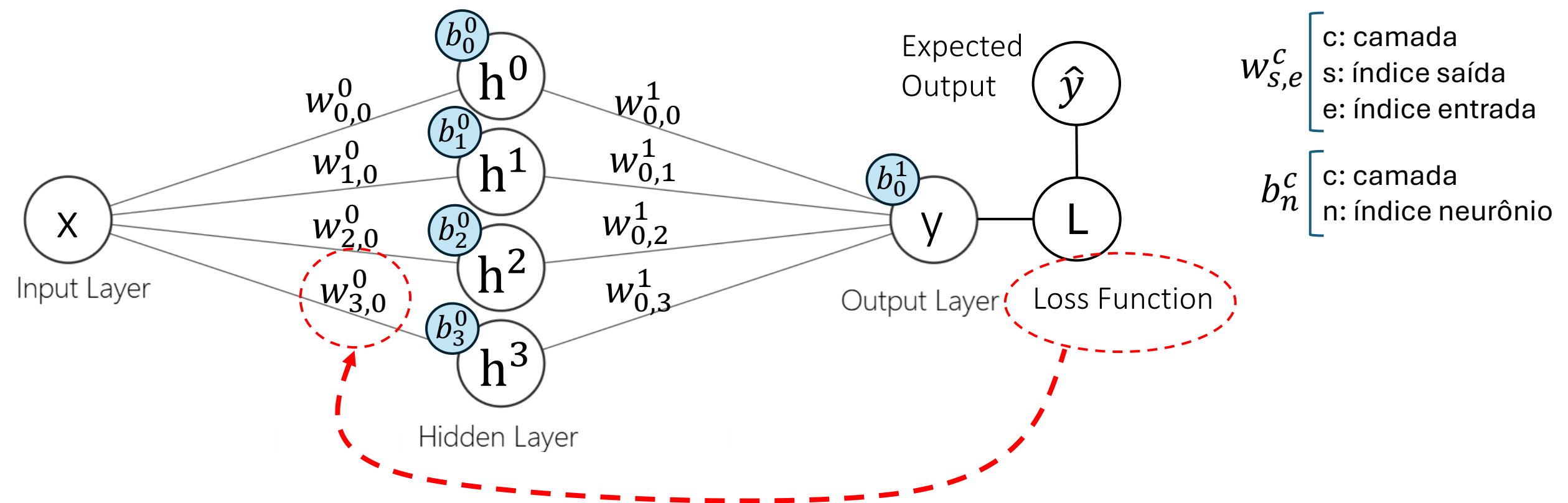
A derivação é análoga para os demais pesos e para o bias, mesmo com camadas intermediárias com mais neurônios.

$$\frac{\partial}{\partial w_{0,k}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + \dots + h^k w_{0,k}^1 + \dots + h^m w_{0,m}^1 + b_0^1]$$

(apenas para lembrar como era a expressão que levou às derivadas)

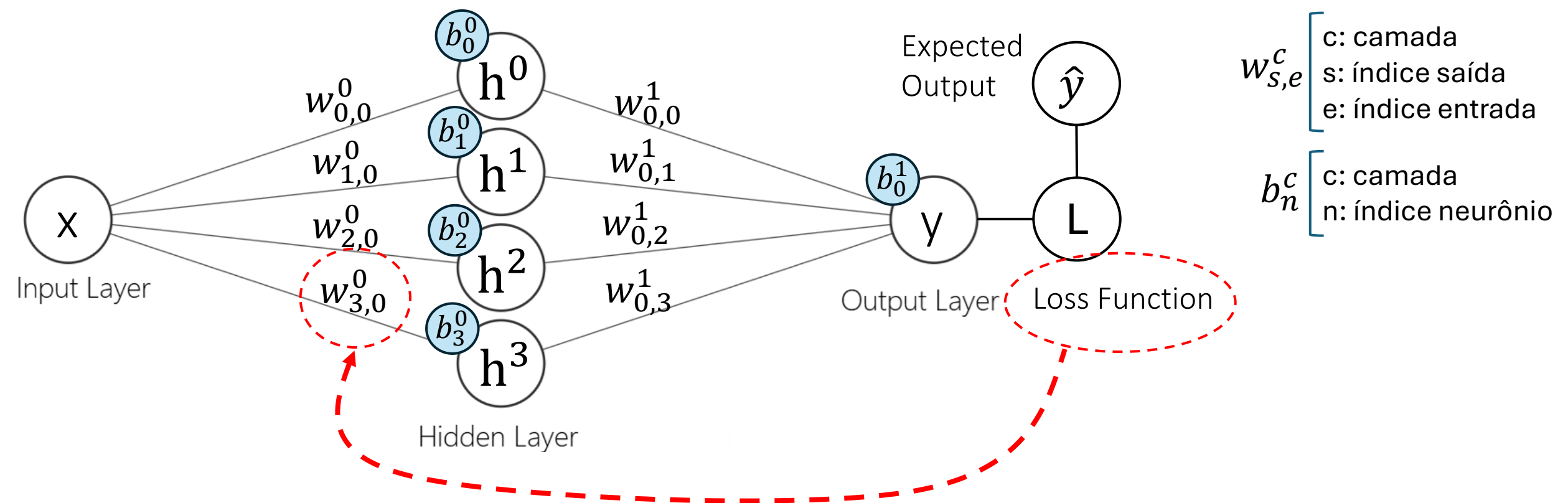


$$\frac{\partial L}{\partial w_{3,0}^0}$$



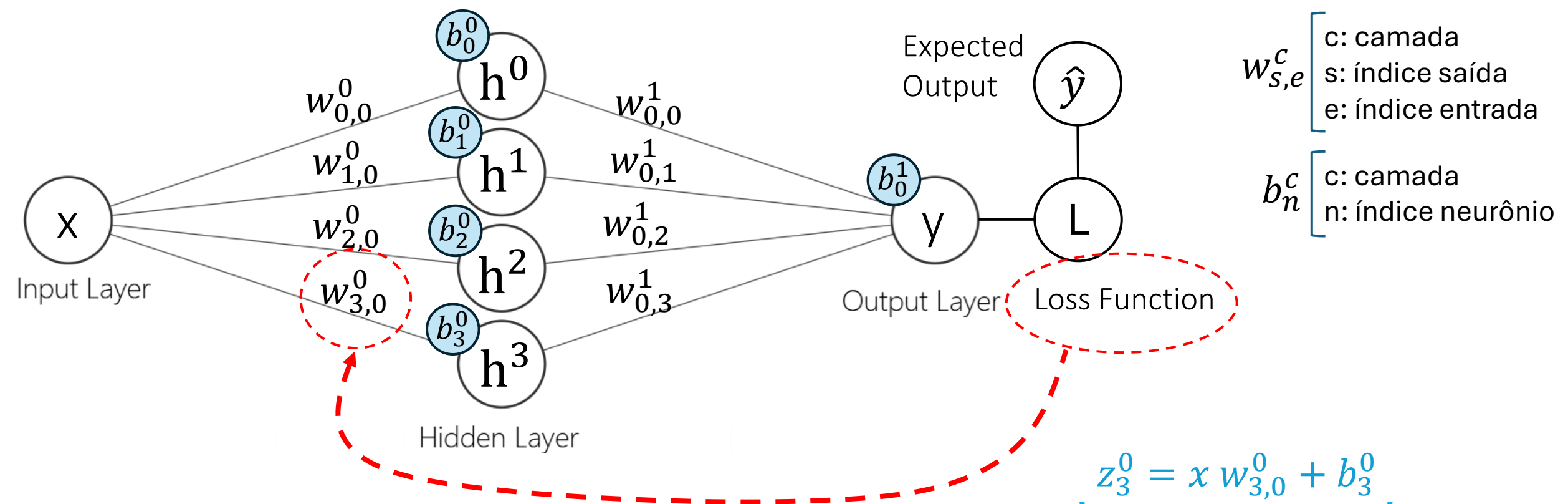
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = \frac{\partial}{\partial w_{3,0}^0} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1]$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

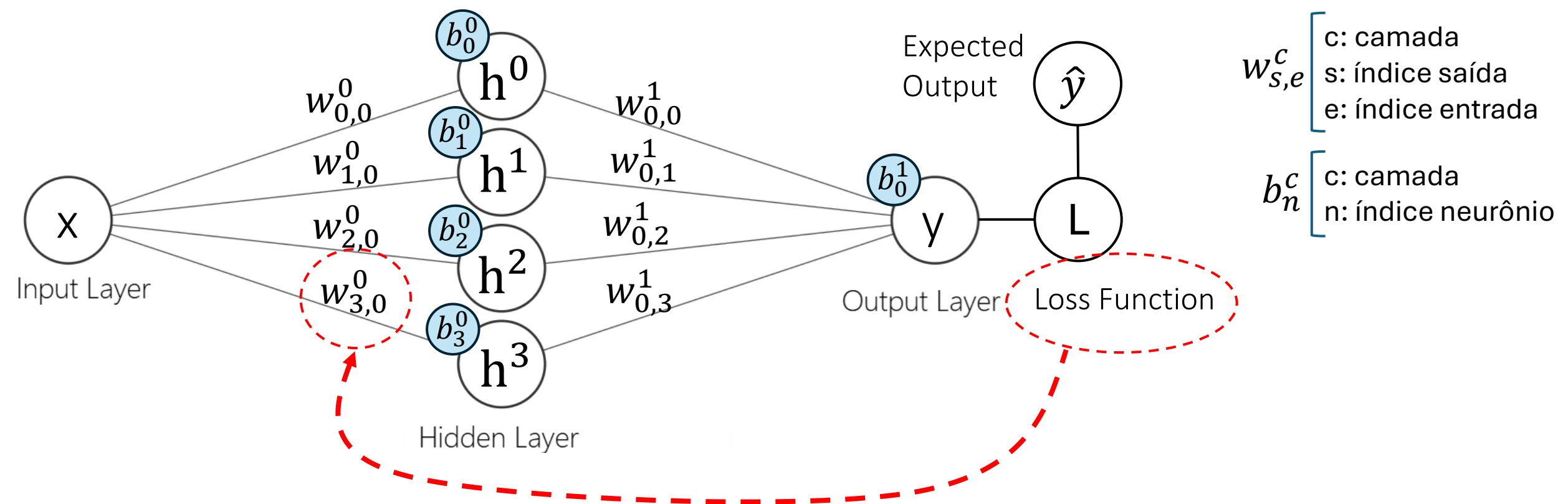
$$\frac{\partial}{\partial w_{3,0}^0} y_i = \frac{\partial}{\partial w_{3,0}^0} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1] = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} h^3 = \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0)$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$

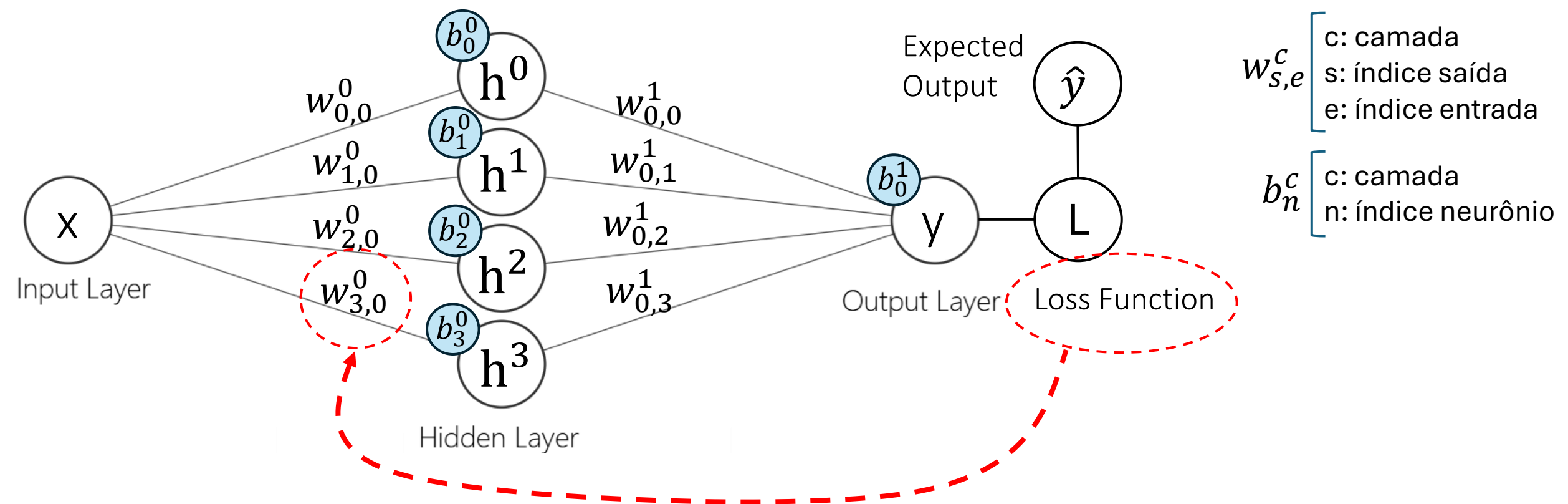


$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$

$$\frac{\partial}{\partial w_{3,0}^0} h^3 = \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0)$$

$$= \sigma'(z_3^0) \frac{\partial}{\partial w_{3,0}^0} z_3^0$$

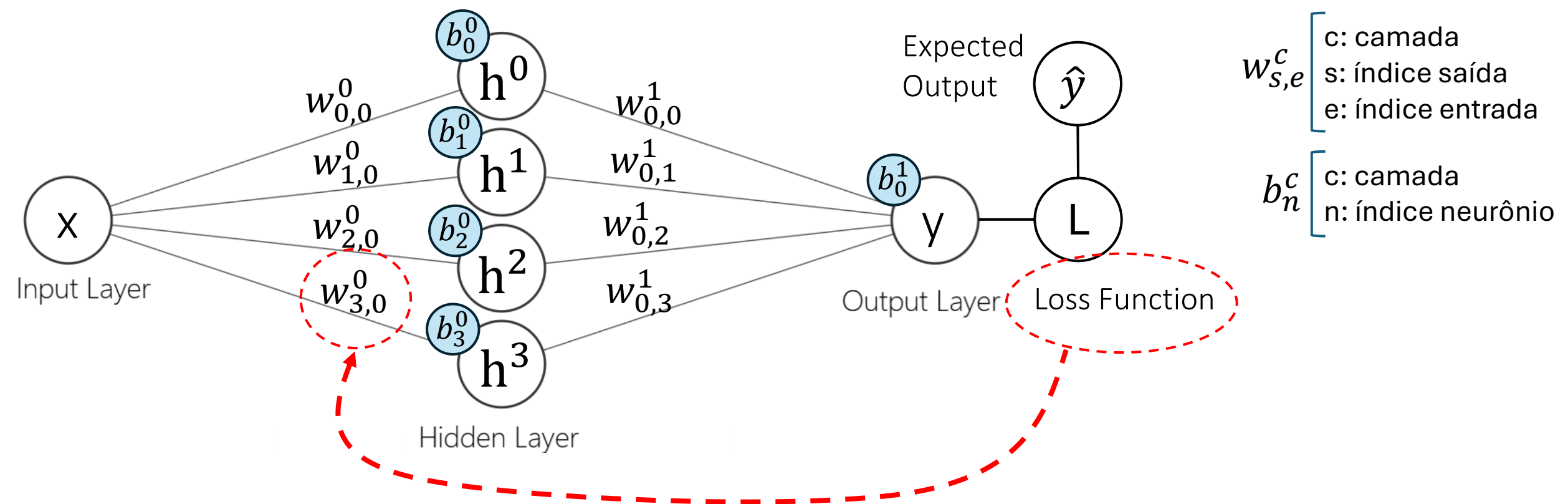


$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} h^3 = \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0)$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$

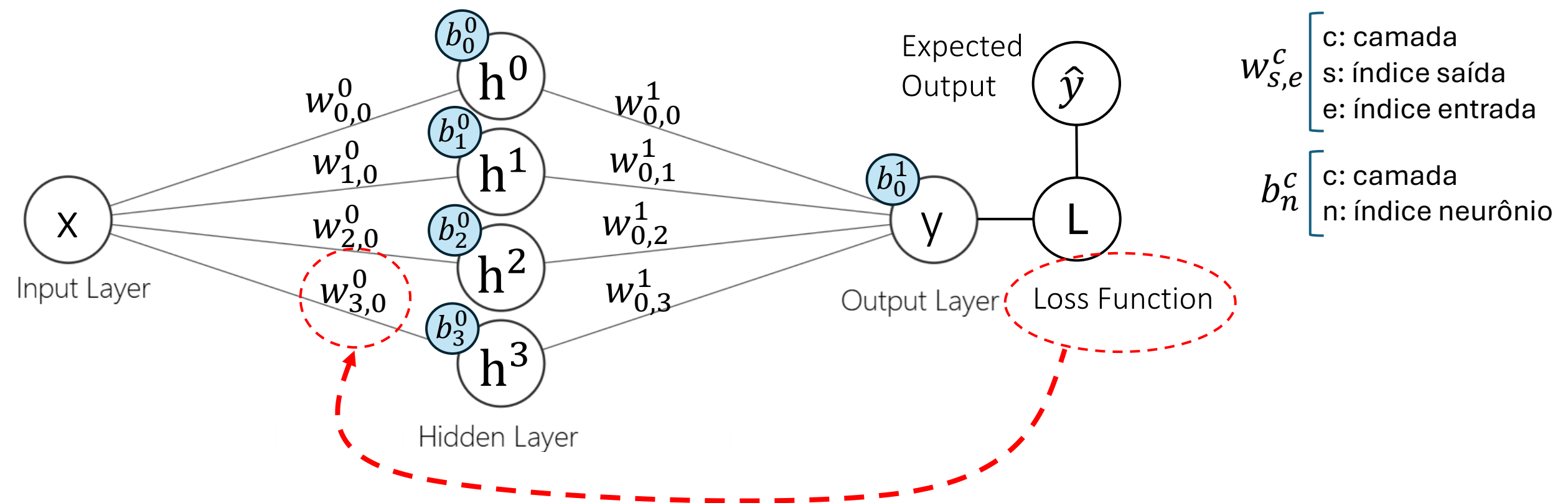
$$= \sigma'(z_3^0) \frac{\partial}{\partial w_{3,0}^0} [x w_{3,0}^0 + b_3^0]$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\begin{aligned} \frac{\partial}{\partial w_{3,0}^0} h^3 &= \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0) \\ &= \sigma'(z_3^0) x \end{aligned}$$

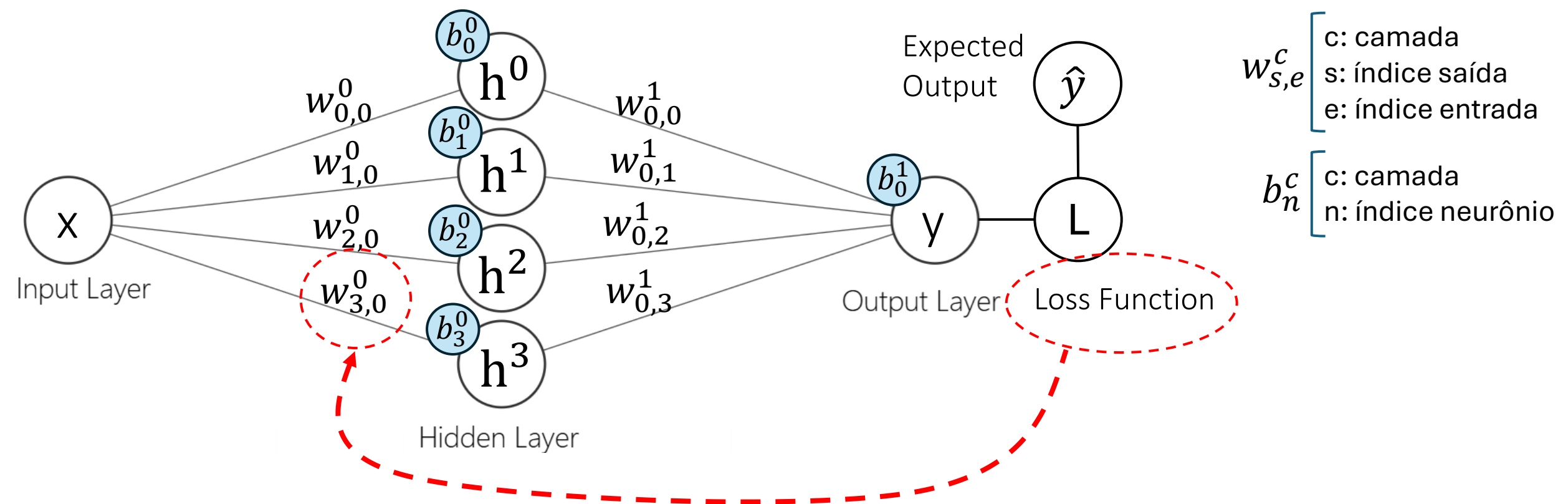
$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \frac{\partial}{\partial w_{3,0}^0} h^3$$

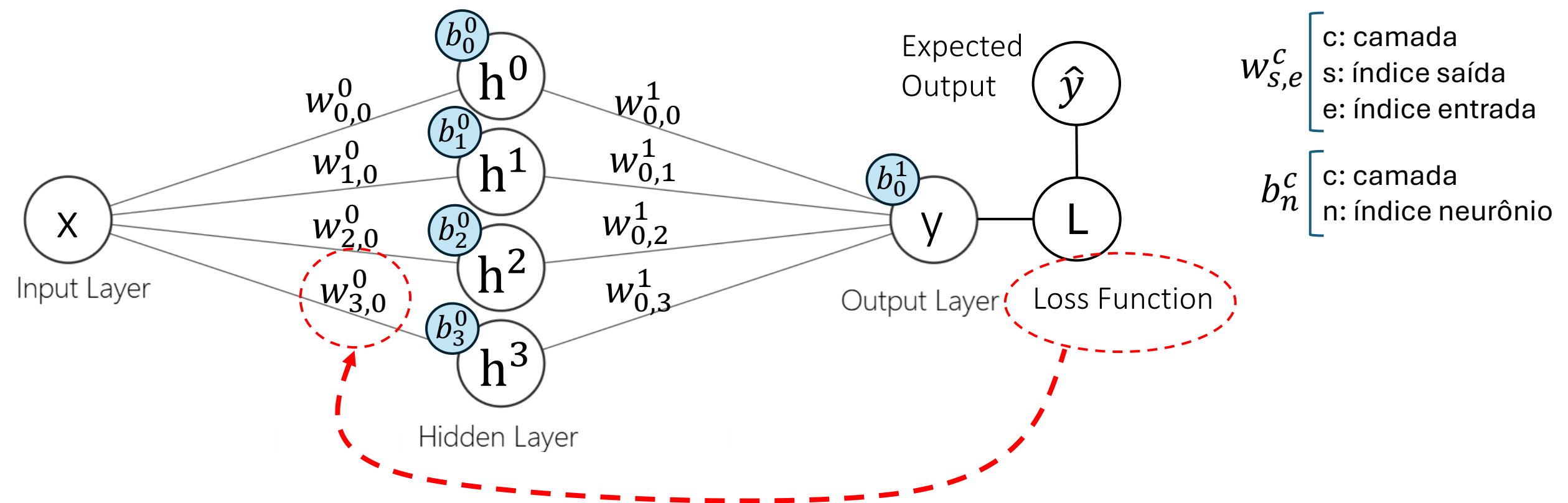
$$\begin{aligned} \frac{\partial}{\partial w_{3,0}^0} h^3 &= \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0) \\ &= \sigma'(z_3^0) x \end{aligned}$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

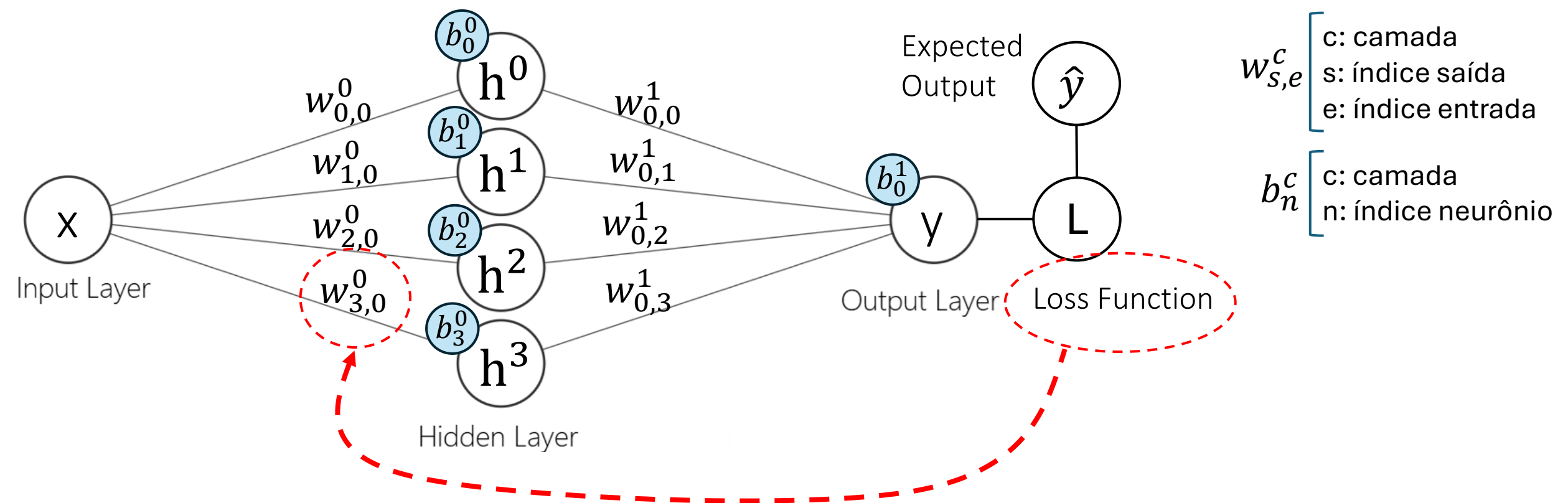
$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \sigma'(z_3^0) x$$

$$\begin{aligned} \frac{\partial}{\partial w_{3,0}^0} h^3 &= \frac{\partial}{\partial w_{3,0}^0} \sigma(z_3^0) \\ &= \sigma'(z_3^0) x \end{aligned}$$



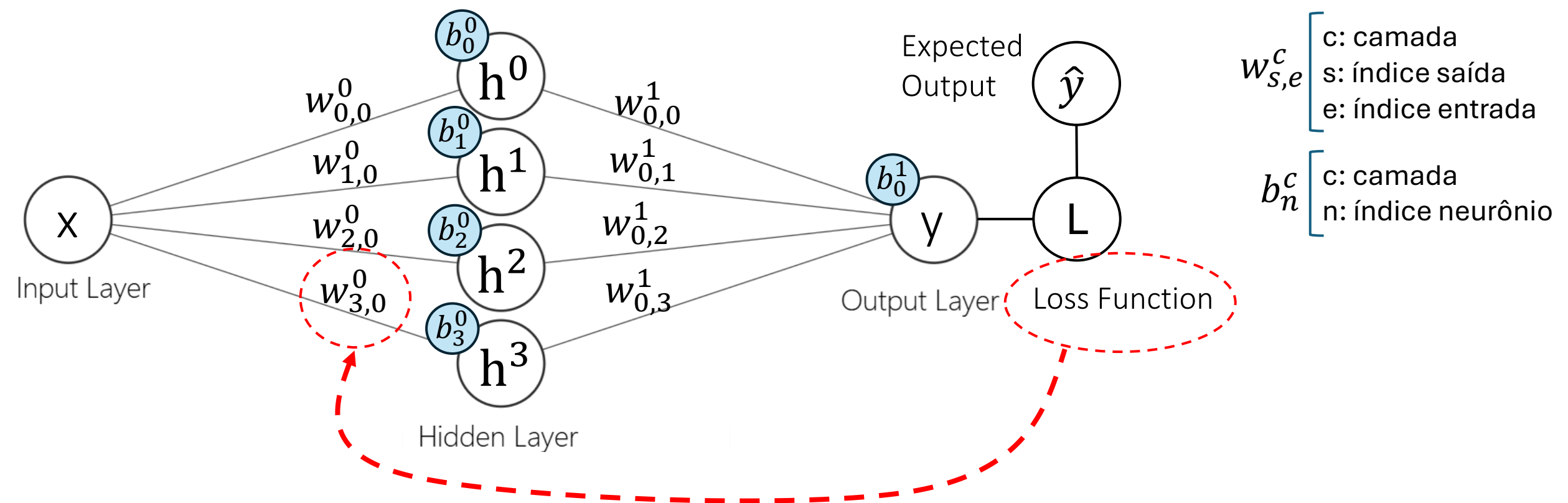
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial}{\partial w_{3,0}^0} y_i$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \sigma'(z_3^0) x$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \sigma'(z_3^0) x$$

$$\frac{\partial}{\partial w_{3,0}^0} y_i = w_{0,3}^1 \sigma'(z_3^0) x$$



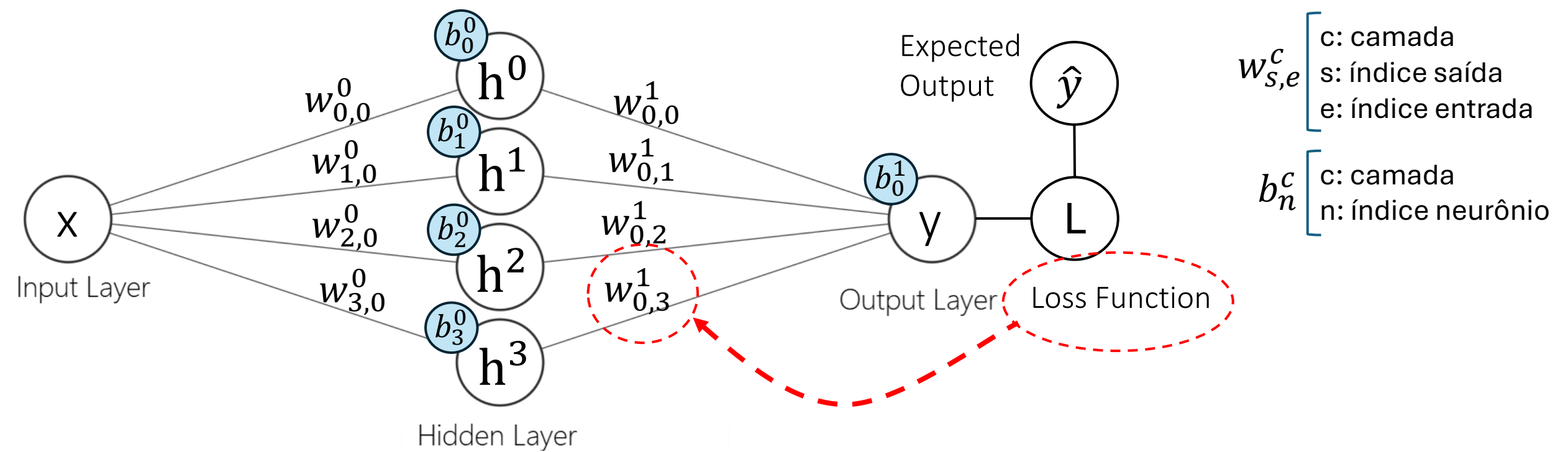
$$\frac{\partial L}{\partial w_{k,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0) x$$

$$\frac{\partial L}{\partial b_k^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0)$$

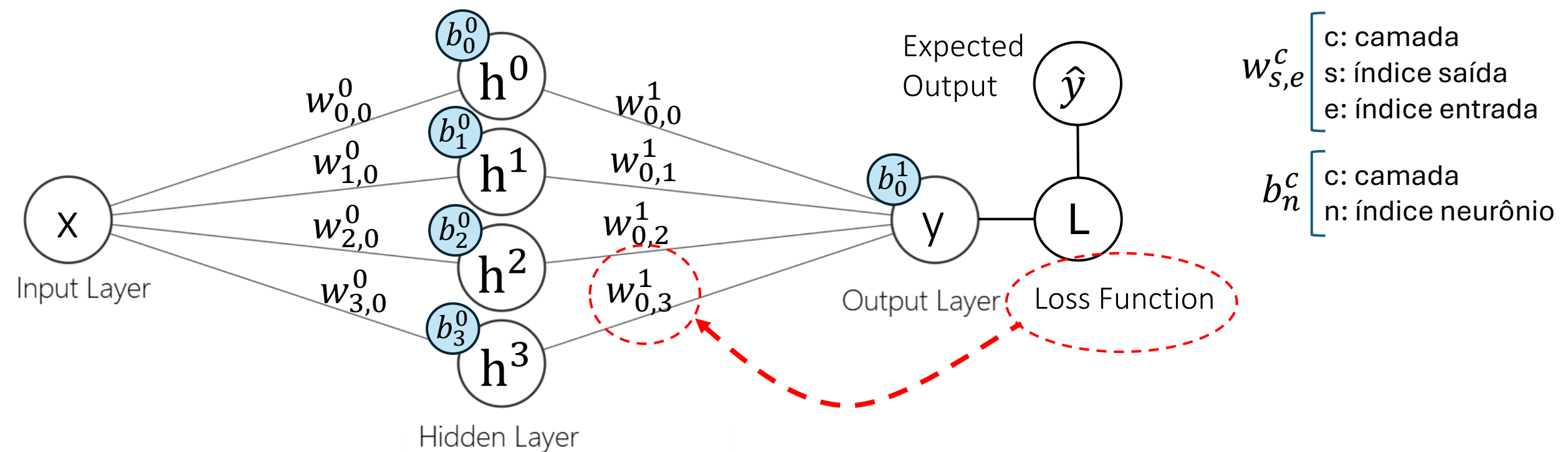
Generalizando para um peso e bias qualquer da camada 0

A próxima seção refaz a derivação anterior resolvendo as derivadas parciais que compõe a regra da cadeia separadamente ao invés de tentar resolver a derivada parcial da função de perda em relação aos pesos diretamente. Este caminho pode ser mais fácil de digerir para algumas pessoas.

Olhando com outra perspectiva: Tornando a regra de cadeia explícita e resolvendo os fatores independentemente



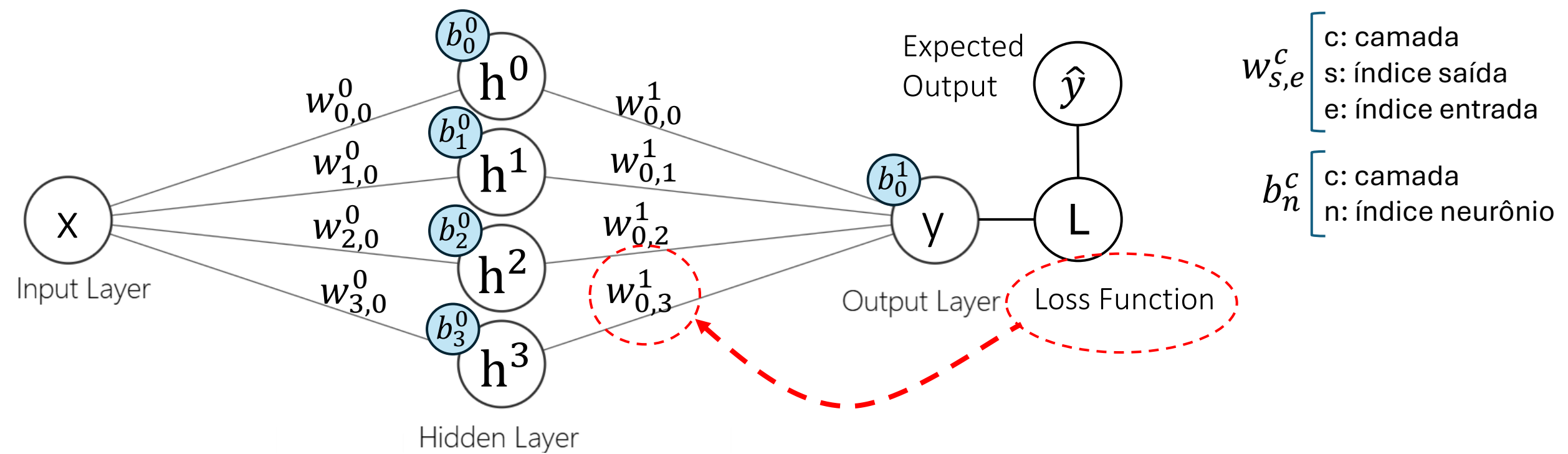
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial w_{0,3}^1}$$



$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial w_{0,3}^1}$$

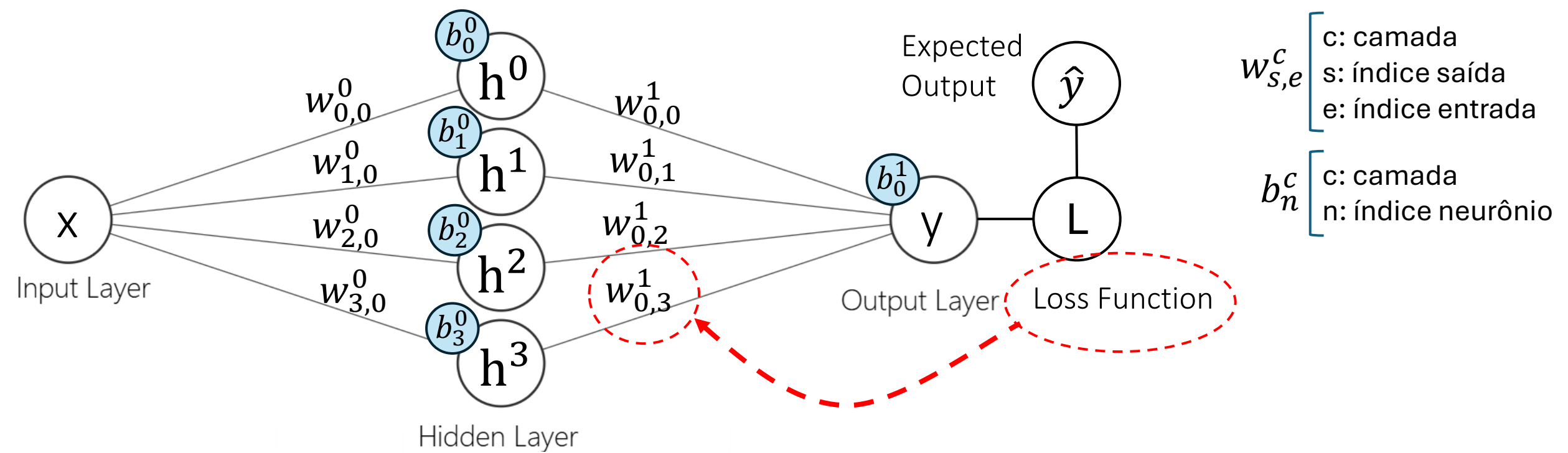
$$\frac{\partial L}{\partial y_i} = \frac{\partial}{\partial y_i} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2$$

Resolvemos a derivada parcial assumindo que y_i é nossa variável de interesse (sem considerar que na verdade queremos a derivada parcial em relação a $w_{0,3}^1$ e que y_i é uma função de $w_{0,3}^1$).



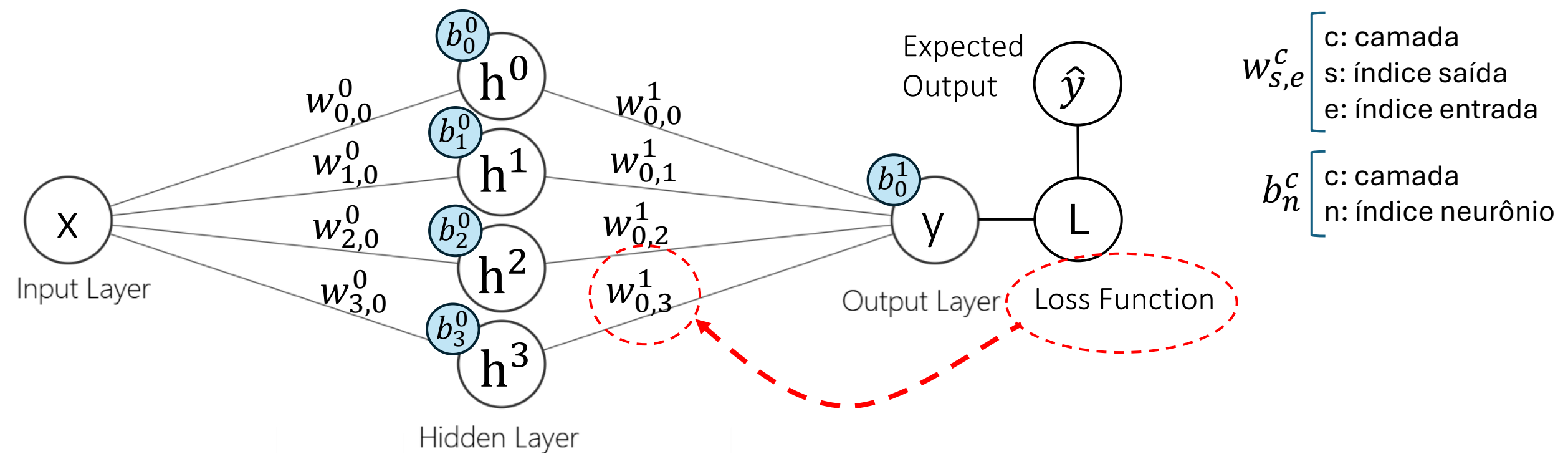
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial}{\partial y_i} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial y_i} (y_i - \hat{y}_i)^2$$



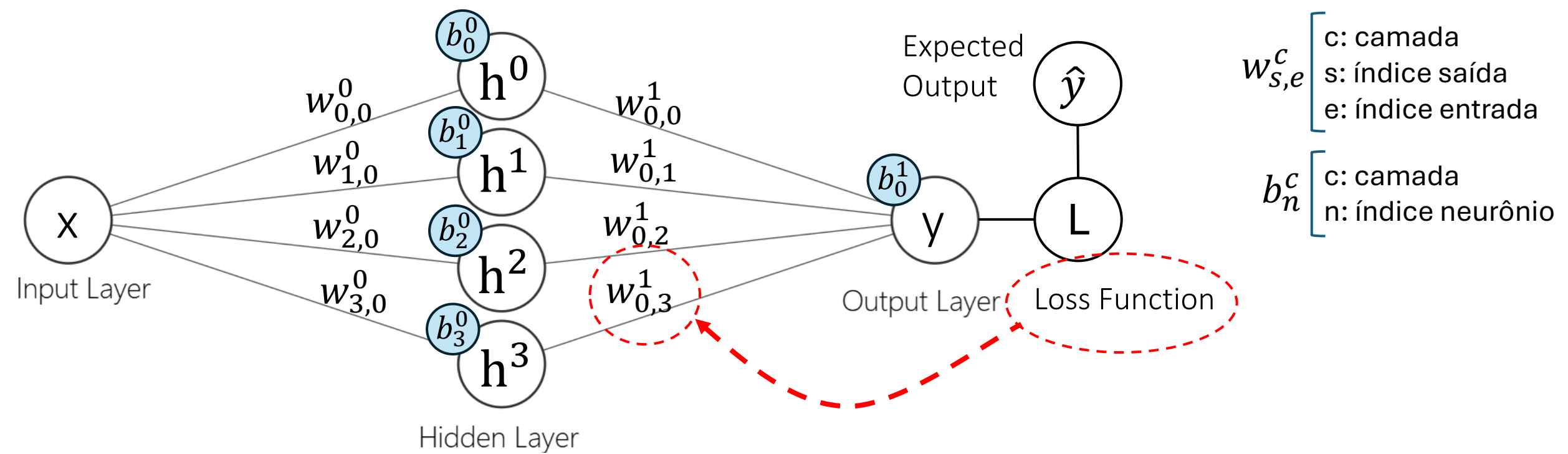
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial}{\partial y_i} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial y_i} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$



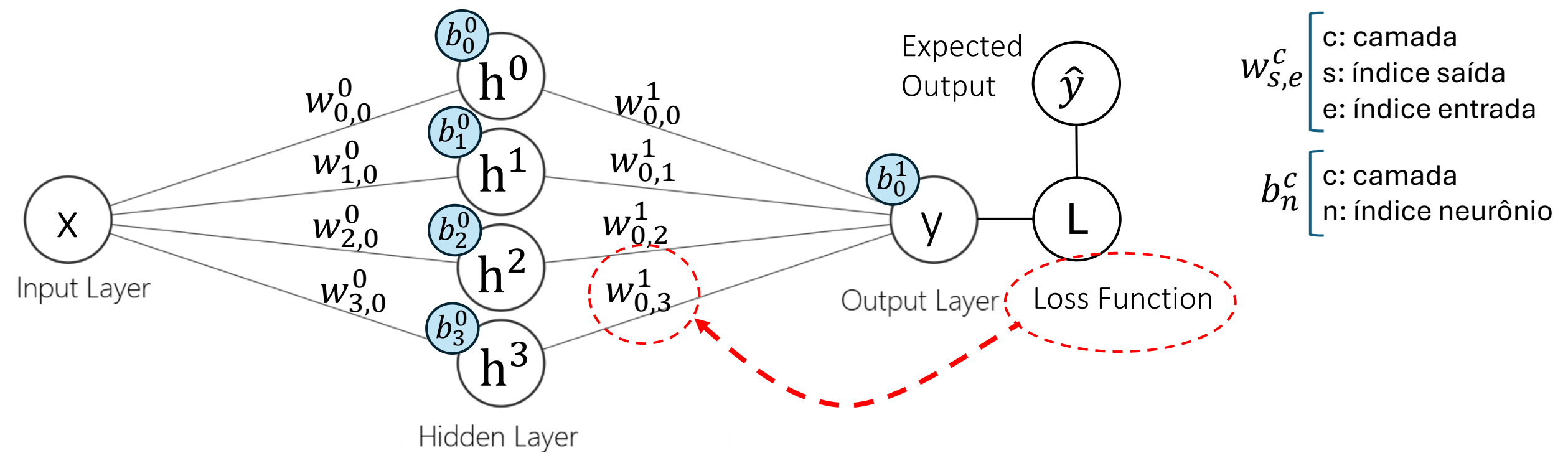
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial}{\partial y_i} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial y_i} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$



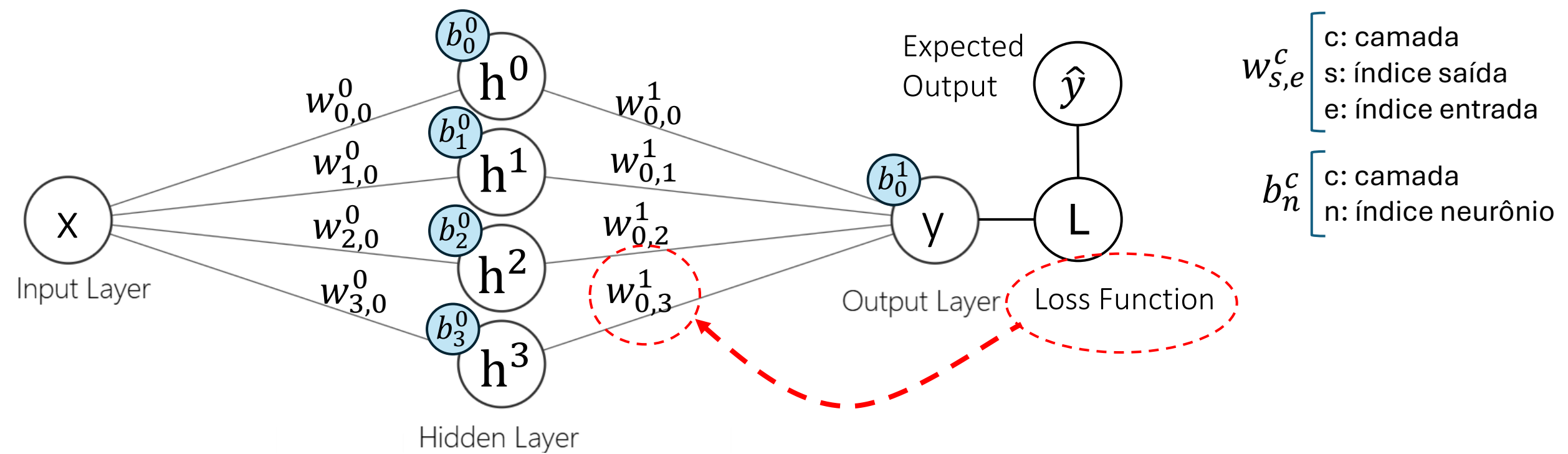
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial L}{\partial y_i} = \frac{\partial}{\partial y_i} \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N \frac{\partial}{\partial y_i} (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$



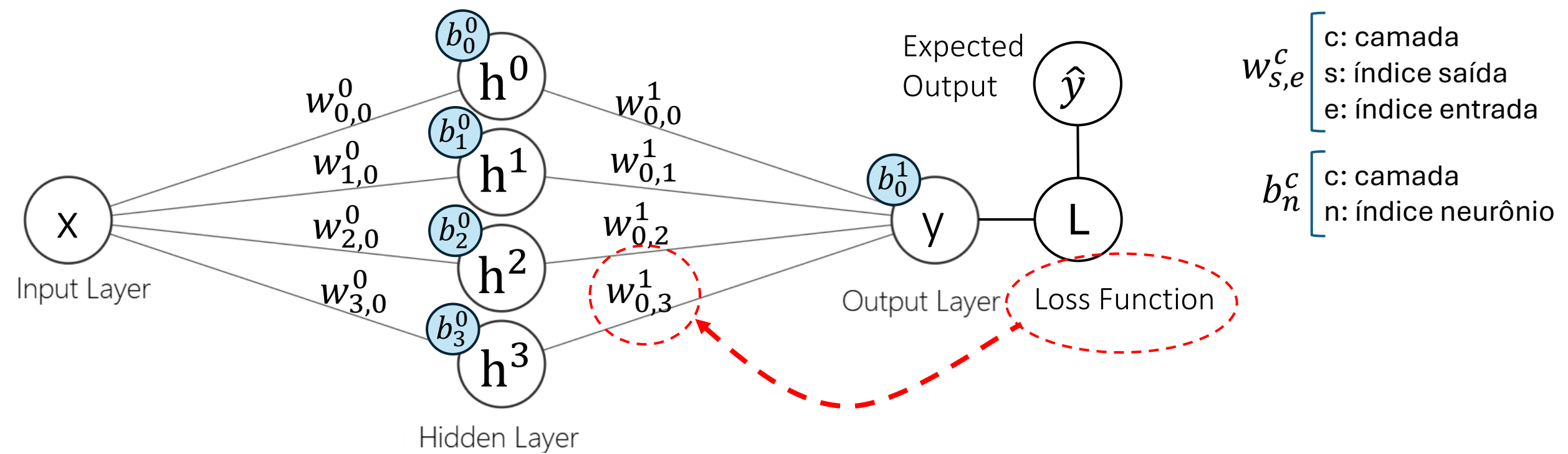
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial y_i}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1]$$



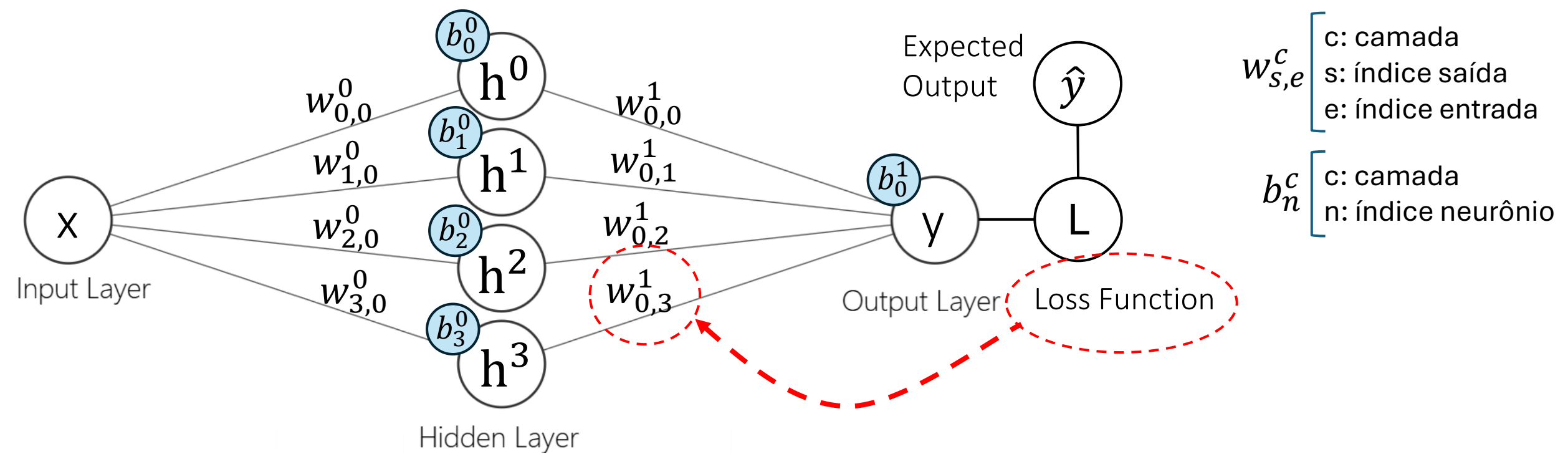
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial y_i}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1] = h^3$$



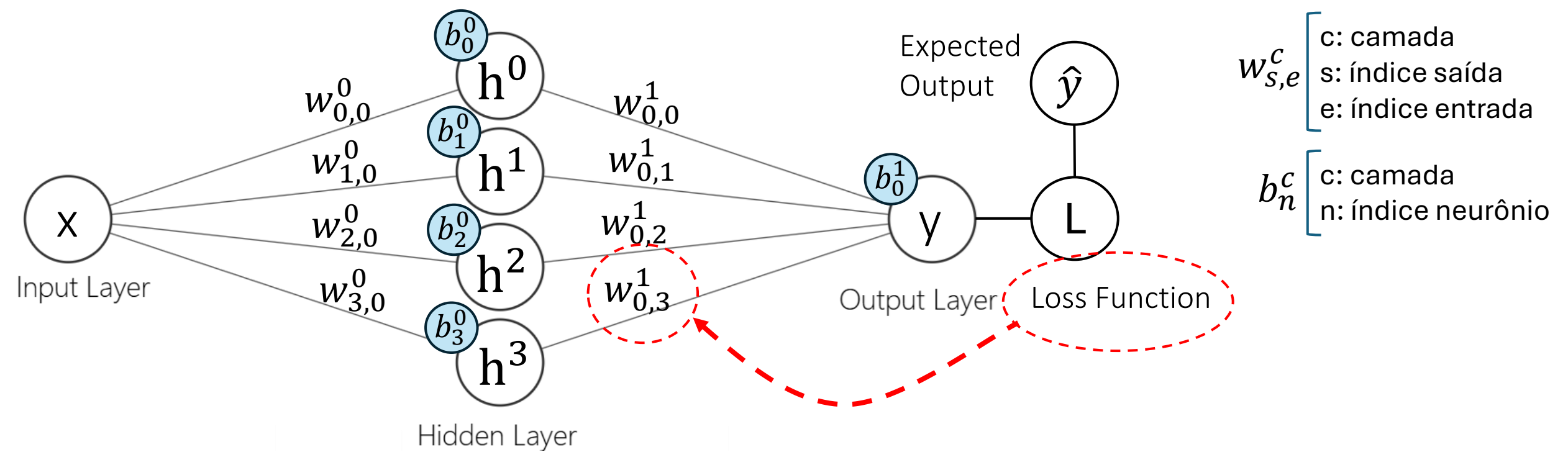
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y_i}{\partial w_{0,3}^1}$$

$$\frac{\partial y_i}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1] = h^3$$

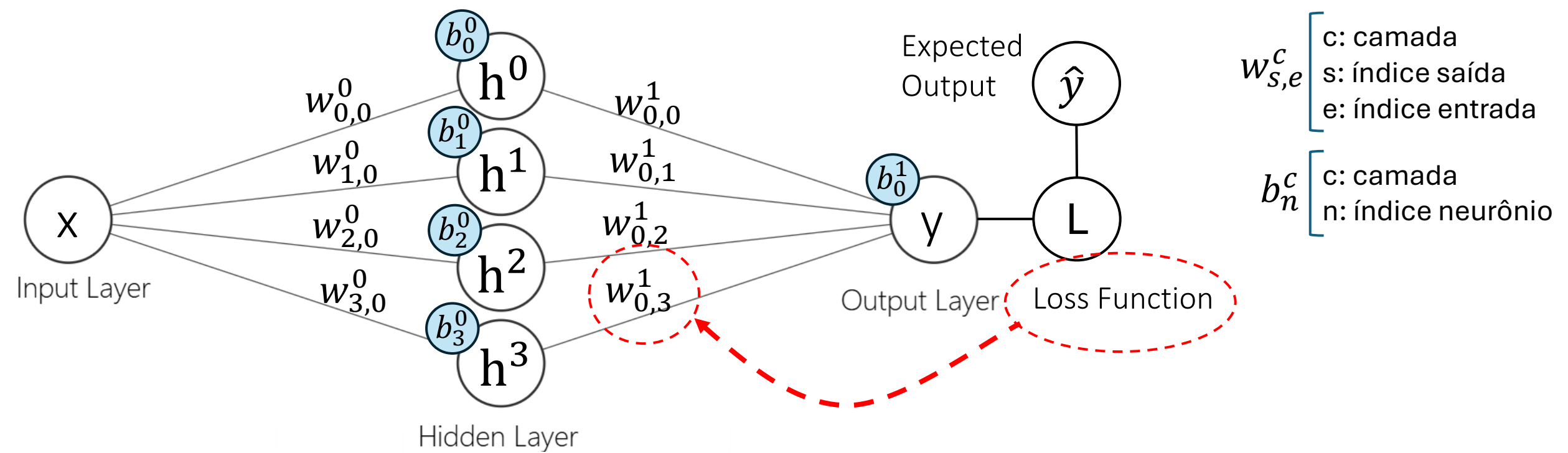


$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$

$$\frac{\partial y_i}{\partial w_{0,3}^1} = \frac{\partial}{\partial w_{0,3}^1} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1] = h^3$$



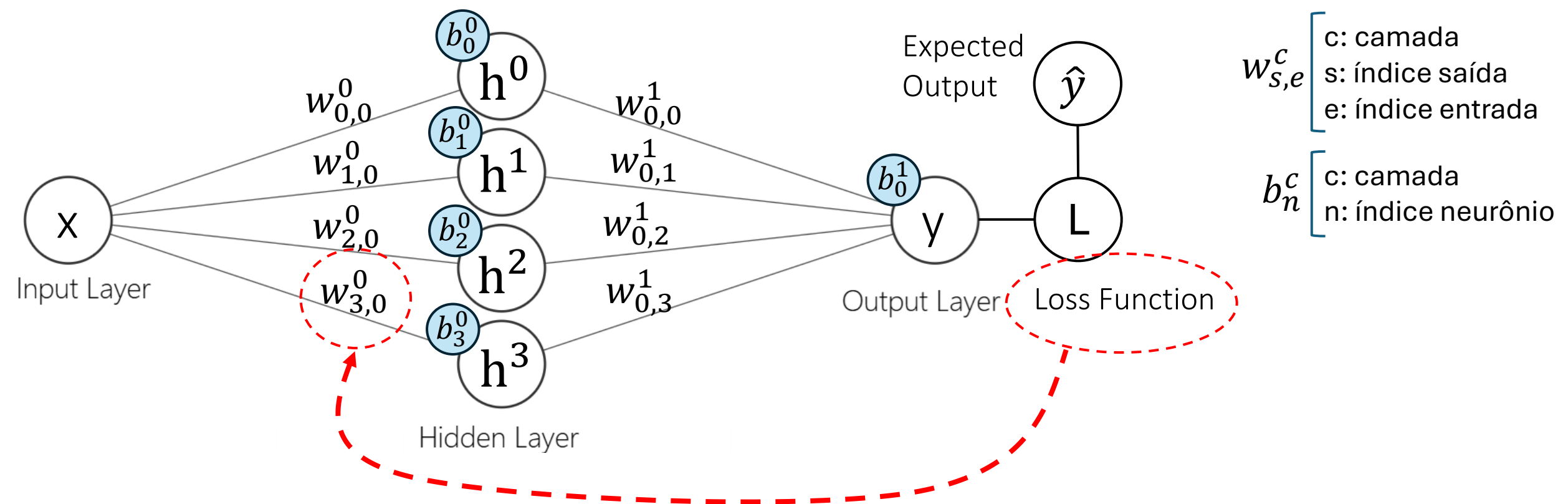
$$\frac{\partial L}{\partial w_{0,3}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^3$$



Generalizando para o k -ésimo peso e para o bias da última camada:

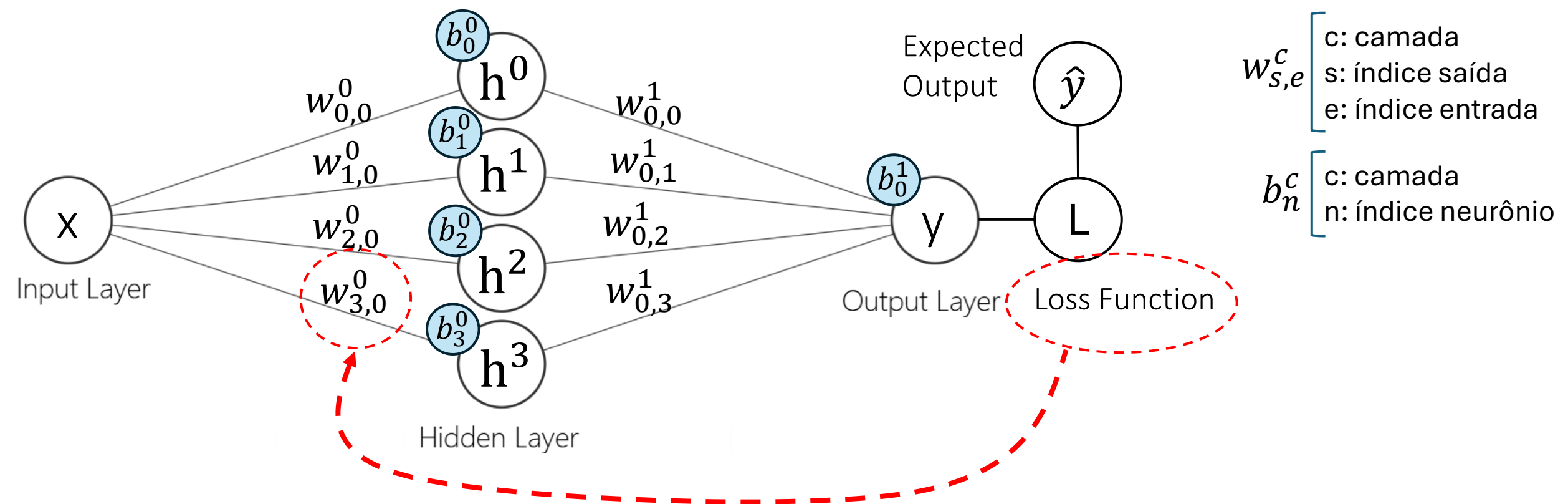
$$\frac{\partial L}{\partial w^1_{0,k}} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^k$$

$$\frac{\partial L}{\partial b^1_0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$



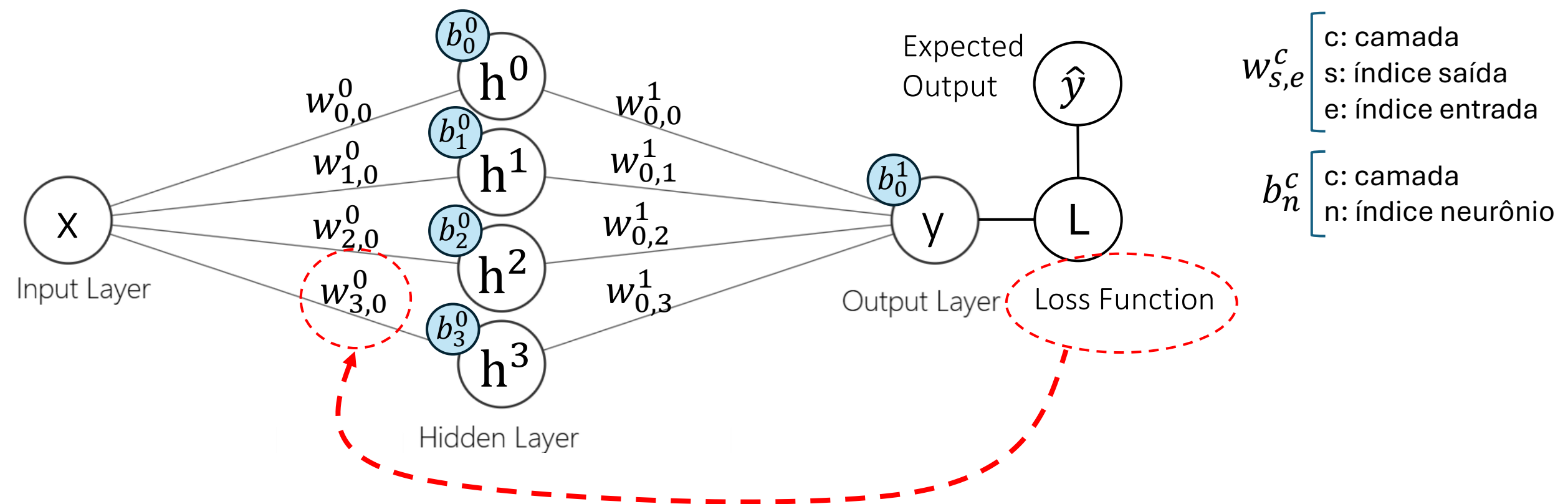
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h^3} \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial L}{\partial y_i} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \quad \text{da derivação da última camada}$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y}{\partial h^3} \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

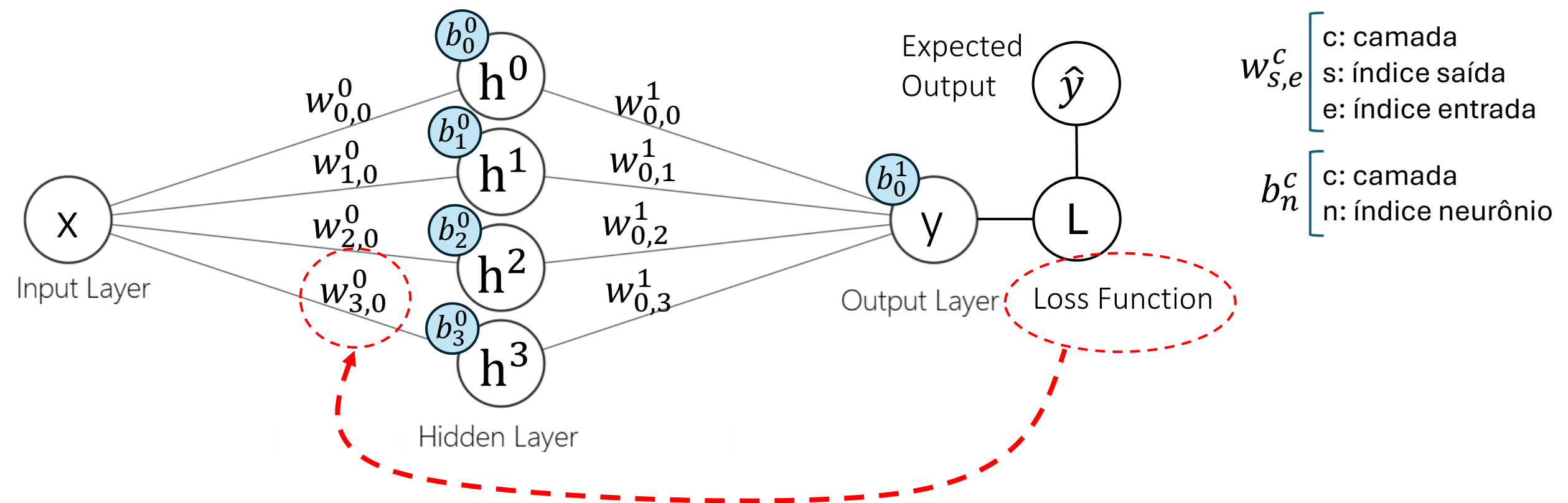
$$\frac{\partial L}{\partial y_i} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \quad \text{da derivação da última camada}$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y}{\partial h^3} \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

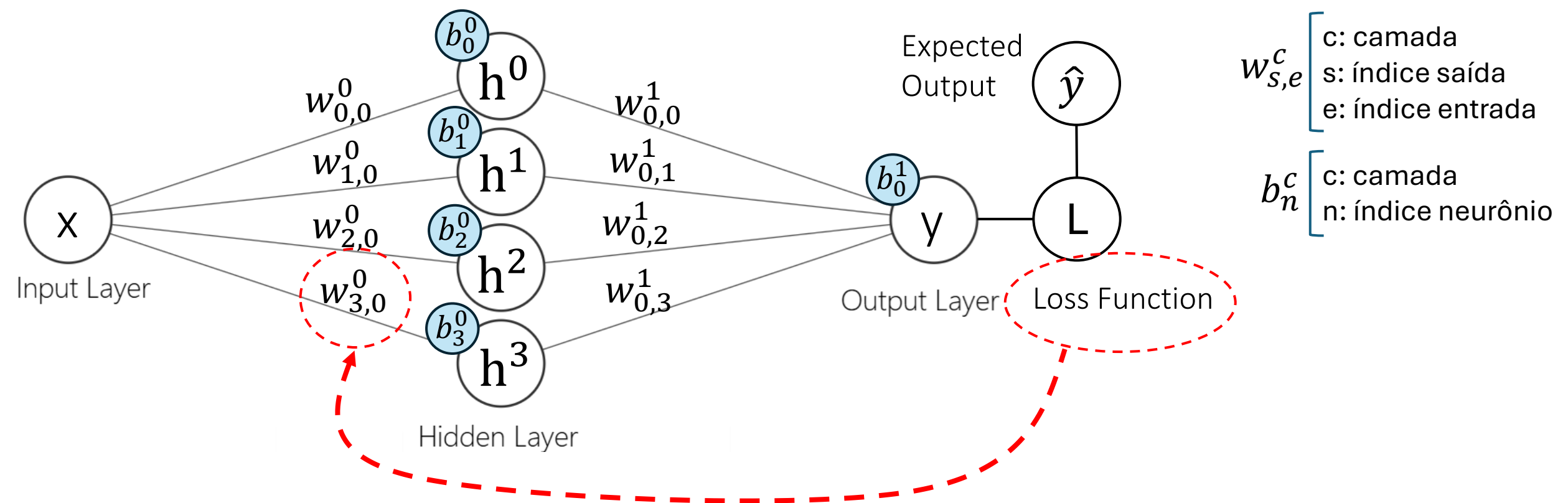
$$\frac{\partial y}{\partial h^3} = \frac{\partial}{\partial h^3} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1]$$

Novamente, resolvemos a derivada parcial como se h^3 fosse nossa variável de interesse



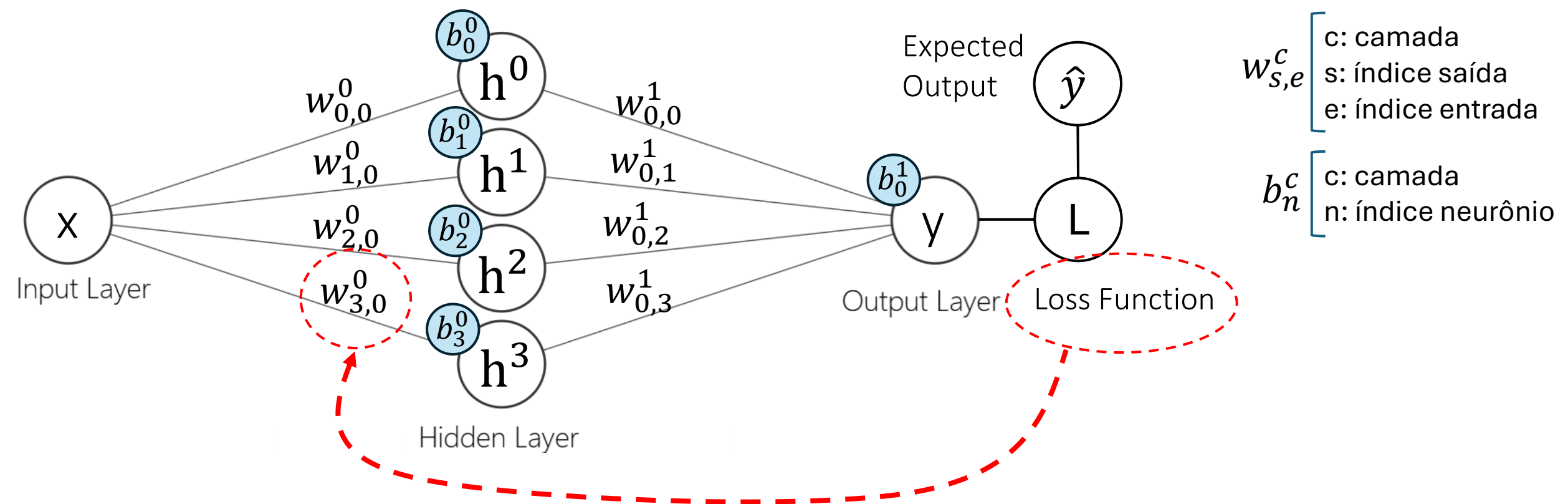
$$\frac{\partial L}{\partial w^0_{3,0}} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) \frac{\partial y}{\partial h^3} \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w^0_{3,0}}$$

$$\frac{\partial y}{\partial h^3} = \frac{\partial}{\partial h^3} [h^0 w^1_{0,0} + h^1 w^1_{0,1} + h^2 w^1_{0,2} + h^3 w^1_{0,3} + b^1_0] = w^1_{0,3}$$



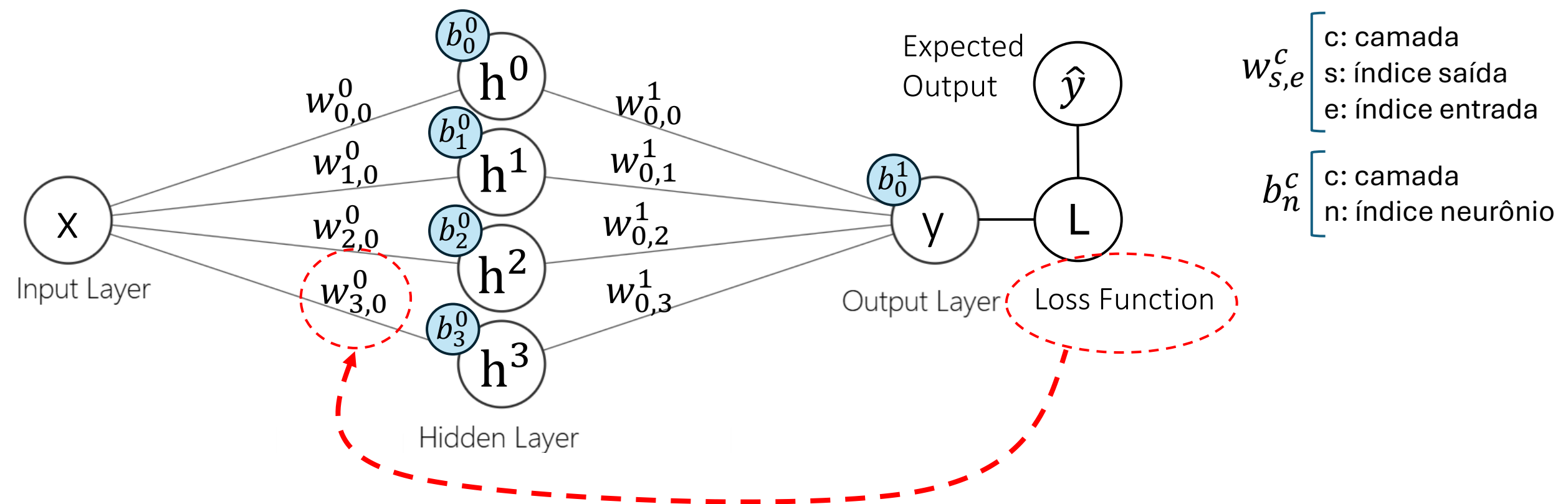
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial y}{\partial h^3} = \frac{\partial}{\partial h^3} [h^0 w_{0,0}^1 + h^1 w_{0,1}^1 + h^2 w_{0,2}^1 + h^3 w_{0,3}^1 + b_0^1] = w_{0,3}^1$$



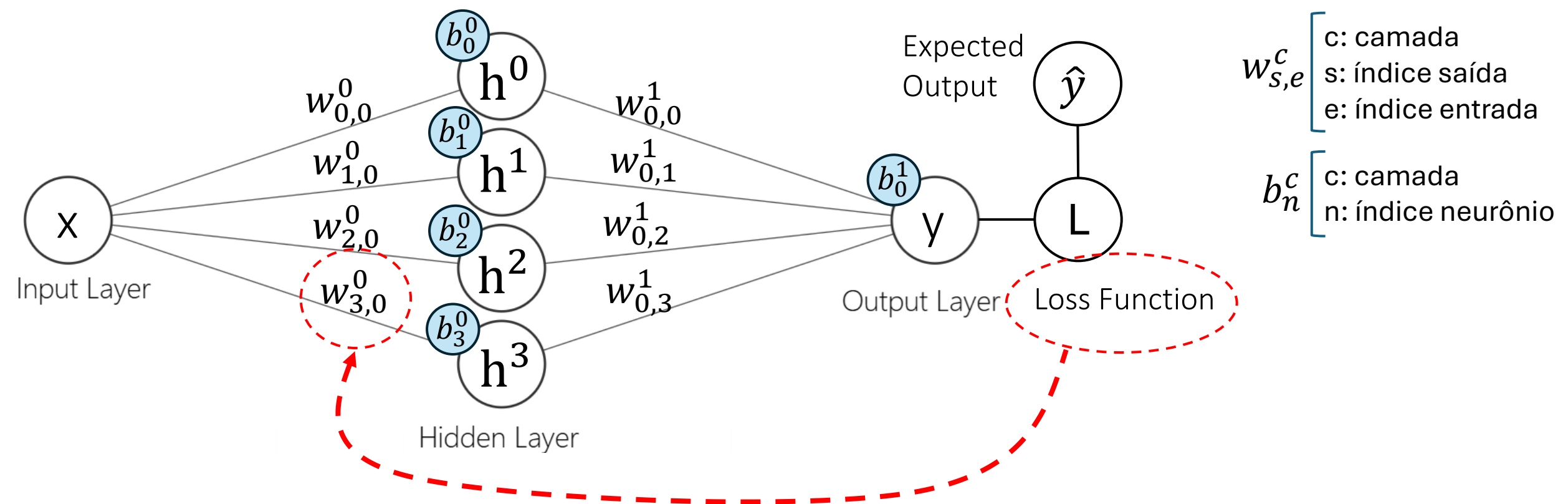
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial h^3}{\partial z^3} = \frac{\partial}{\partial z^3} \sigma(z_3^0)$$



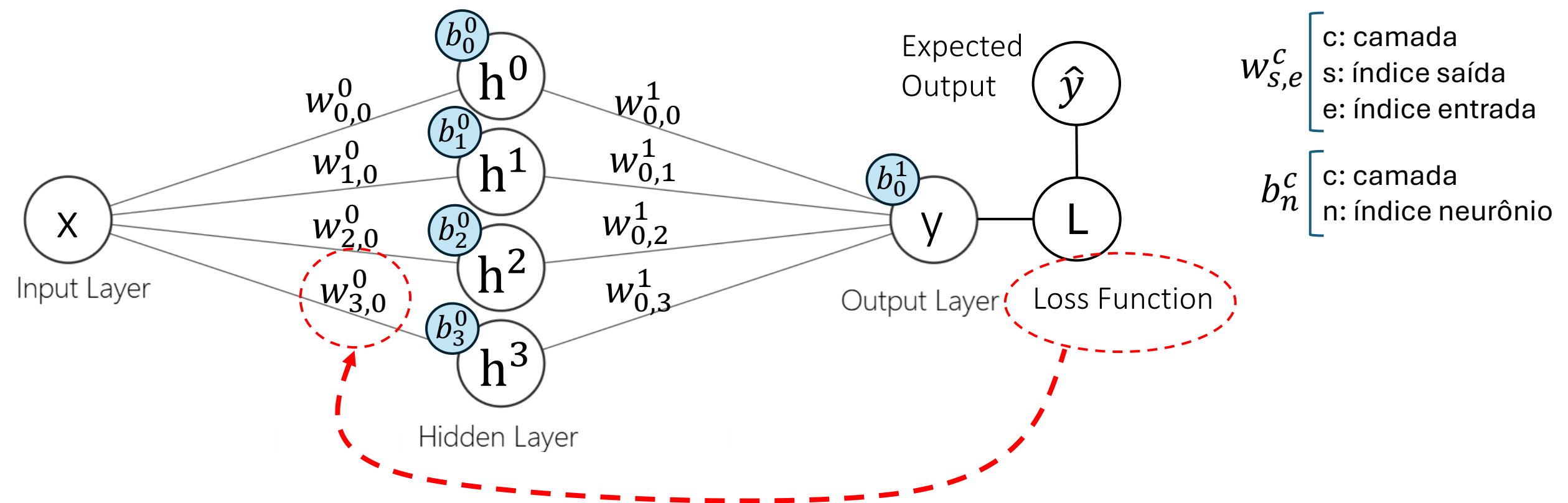
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \frac{\partial h^3}{\partial z^3} \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial h^3}{\partial z^3} = \frac{\partial}{\partial z^3} \sigma(z_3^0) = \sigma'(z_3^0)$$



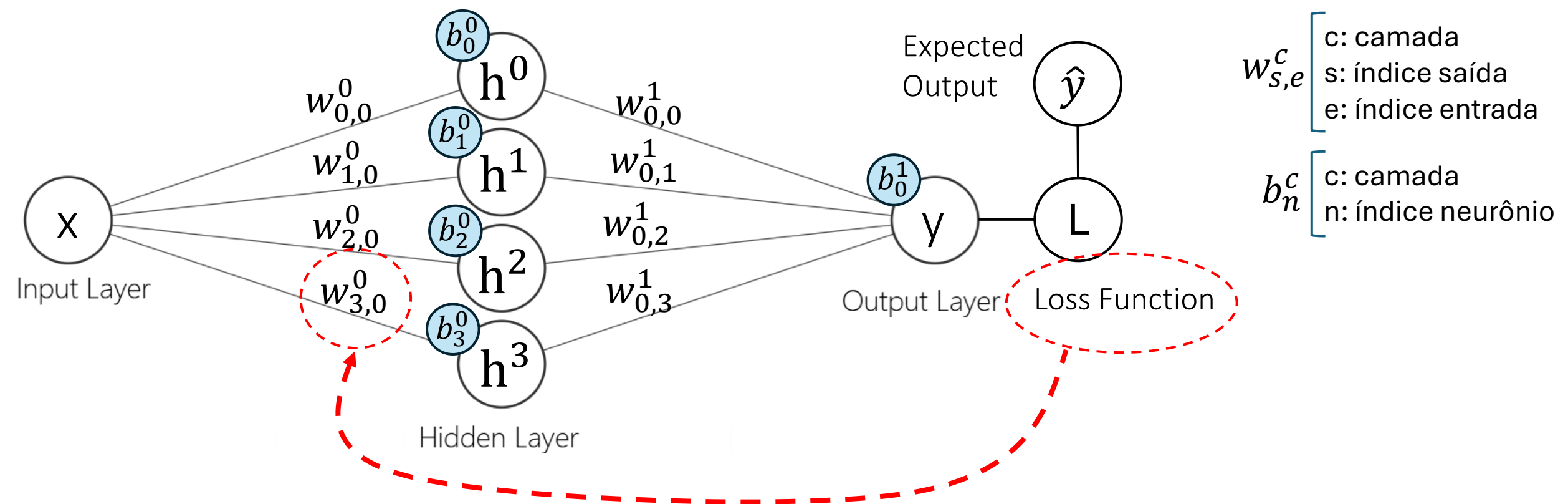
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \sigma'(z_3^0) \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial h^3}{\partial z^3} = \frac{\partial}{\partial z^3} \sigma(z_3^0) = \sigma'(z_3^0)$$



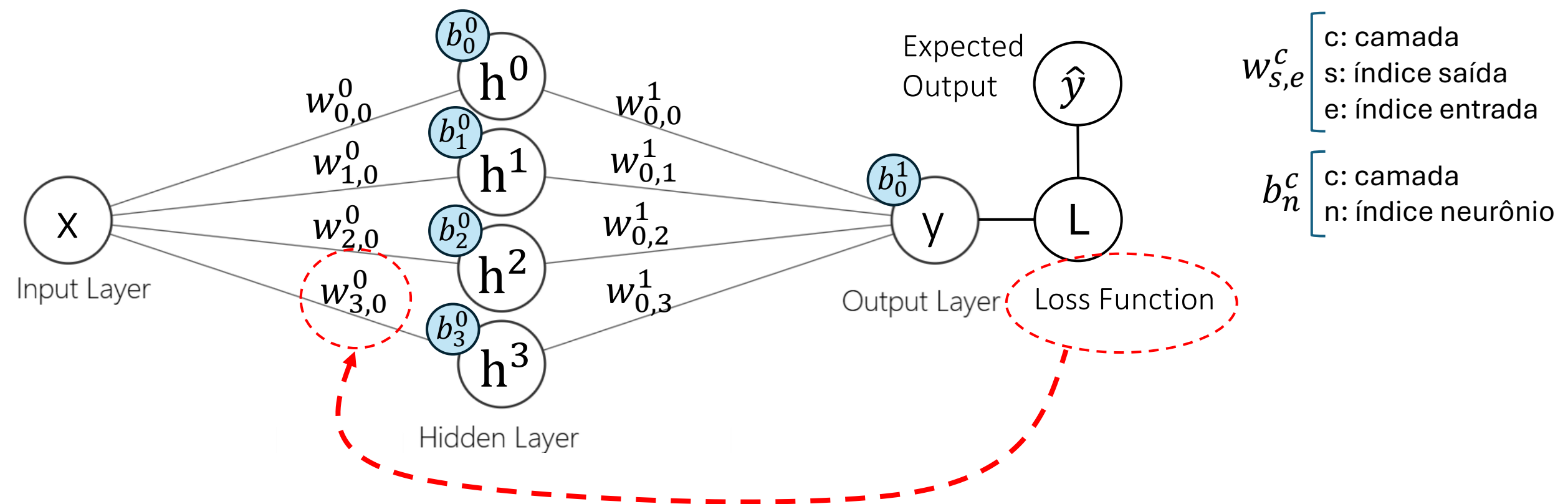
$$\frac{\partial L}{\partial w^0_{3,0}} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w^1_{0,3} \sigma'(z^0_3) \frac{\partial z^3}{\partial w^0_{3,0}}$$

$$\frac{\partial z^3}{\partial w^0_{3,0}} = \frac{\partial}{\partial w^0_{3,0}} [x w^0_{3,0} + b^0_3]$$



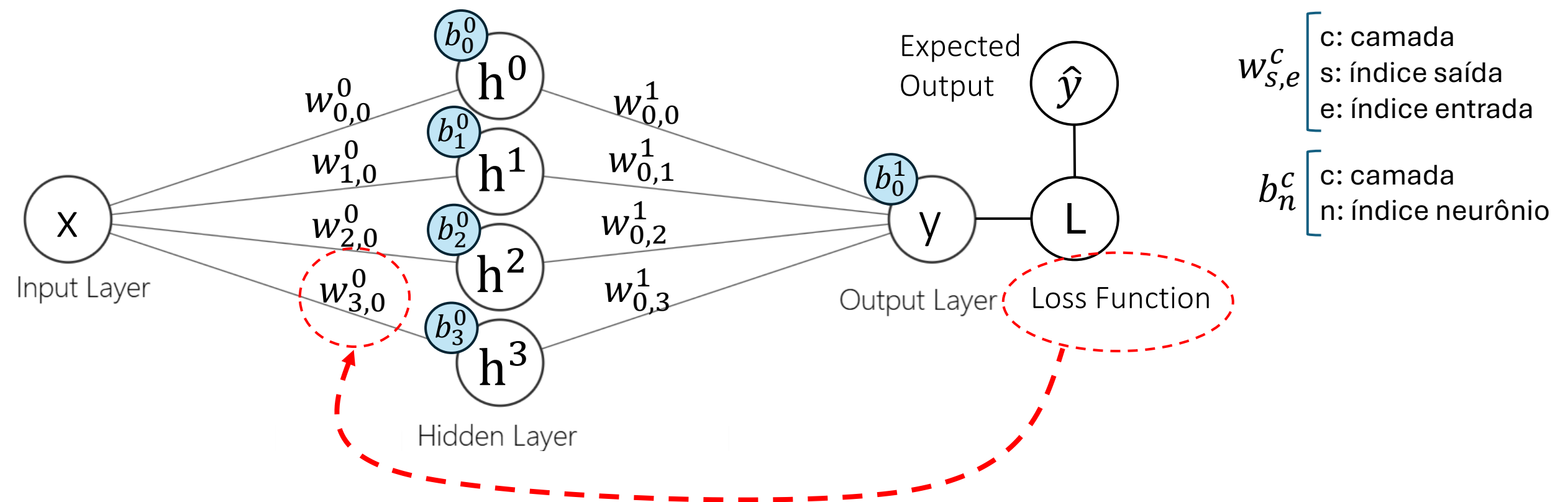
$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \sigma'(z_3^0) \frac{\partial z^3}{\partial w_{3,0}^0}$$

$$\frac{\partial z^3}{\partial w_{3,0}^0} = \frac{\partial}{\partial w_{3,0}^0} [x w_{3,0}^0 + b_3^0] = x$$

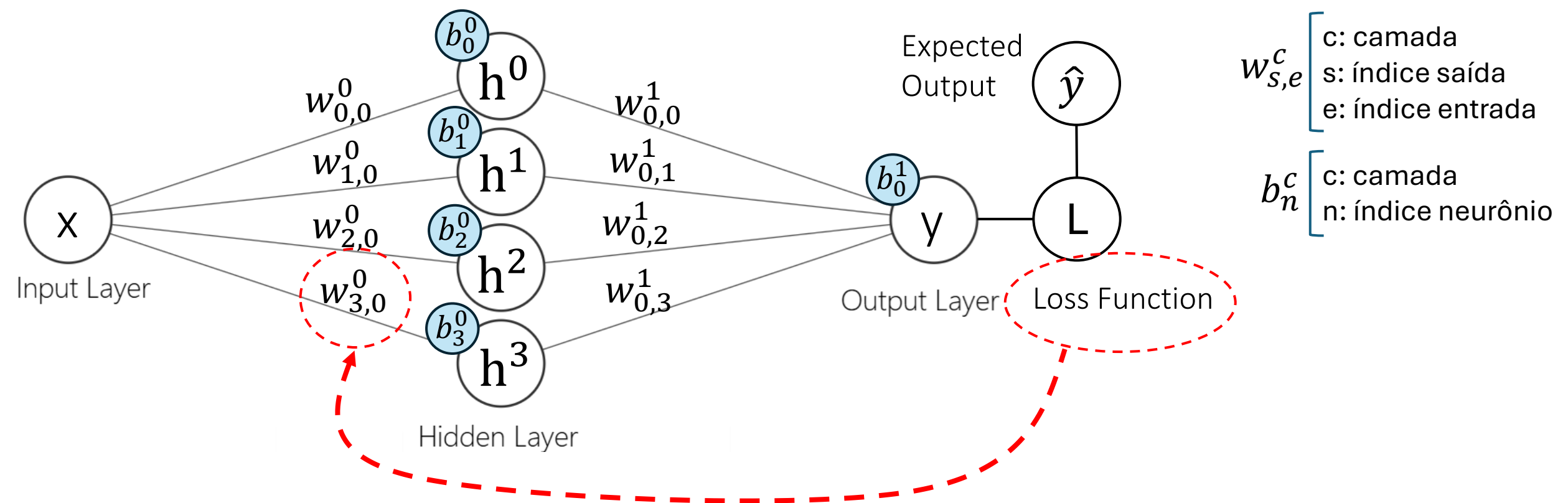


$$\frac{\partial L}{\partial w^0_{3,0}} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w^1_{0,3} \sigma'(z^0_3) x$$

$$\frac{\partial z^3}{\partial w^0_{3,0}} = \frac{\partial}{\partial w^0_{3,0}} [x w^0_{3,0} + b^0_3] = x$$



$$\frac{\partial L}{\partial w_{3,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,3}^1 \sigma'(z_3^0) x$$

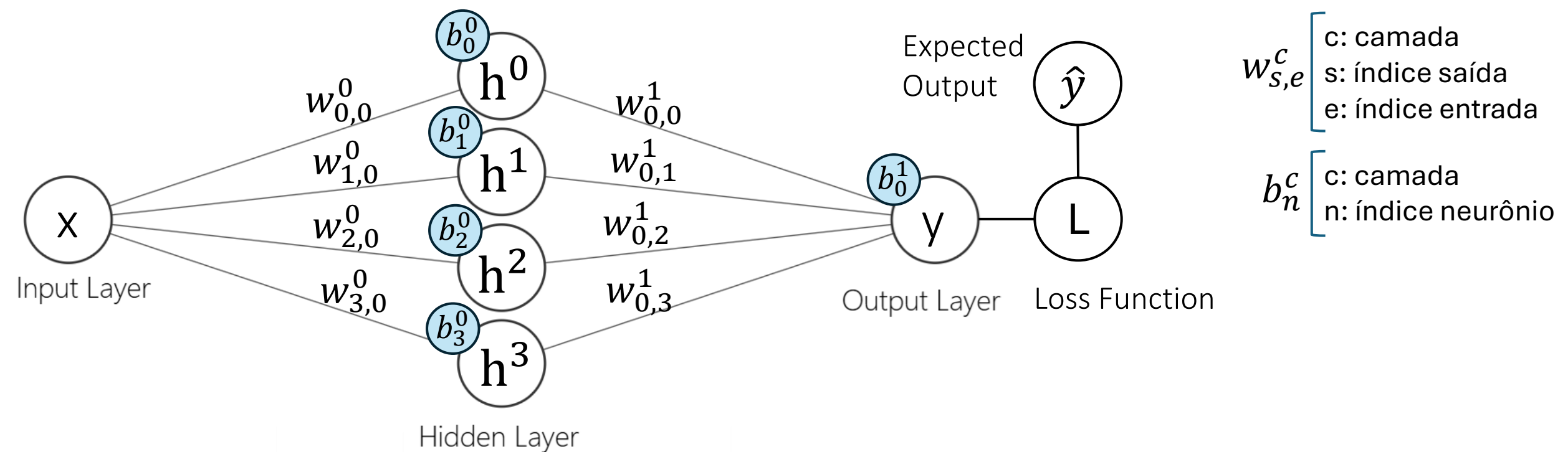


Generalizando para o k -ésimo peso e para o bias da camada oculta:

$$\frac{\partial L}{\partial w_{k,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0) x$$

$$\frac{\partial L}{\partial b_k^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0)$$

Resumo do Resultado Final



Última camada:

$$\frac{\partial L}{\partial w_{0,k}^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) h^k$$

$$\frac{\partial L}{\partial b_0^1} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i)$$

Camada oculta:

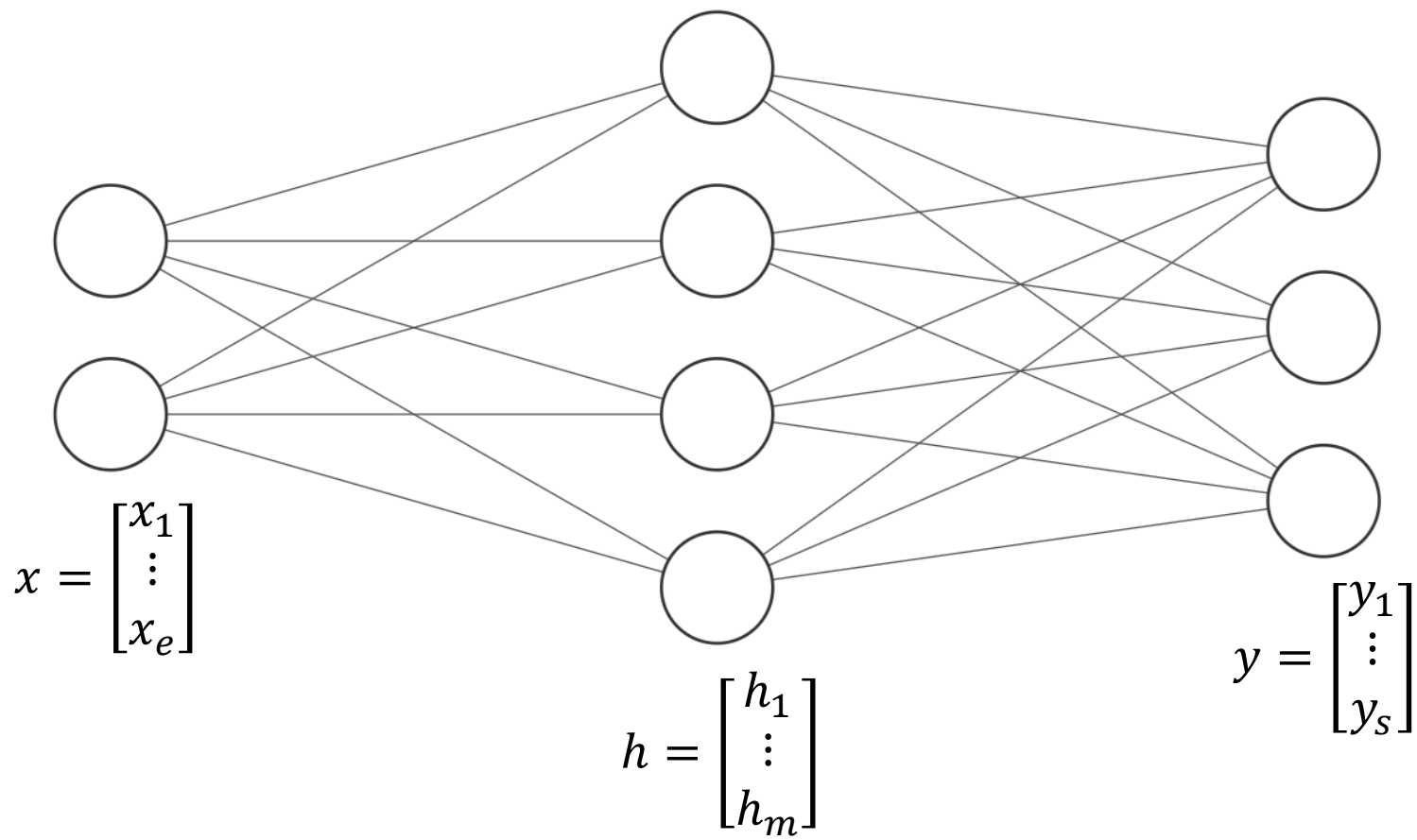
$$\frac{\partial L}{\partial w_{k,0}^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0) x$$

$$\frac{\partial L}{\partial b_k^0} = \frac{1}{N} \sum_{i=0}^N 2(y_i - \hat{y}_i) w_{0,k}^1 \sigma'(z_k^0)$$

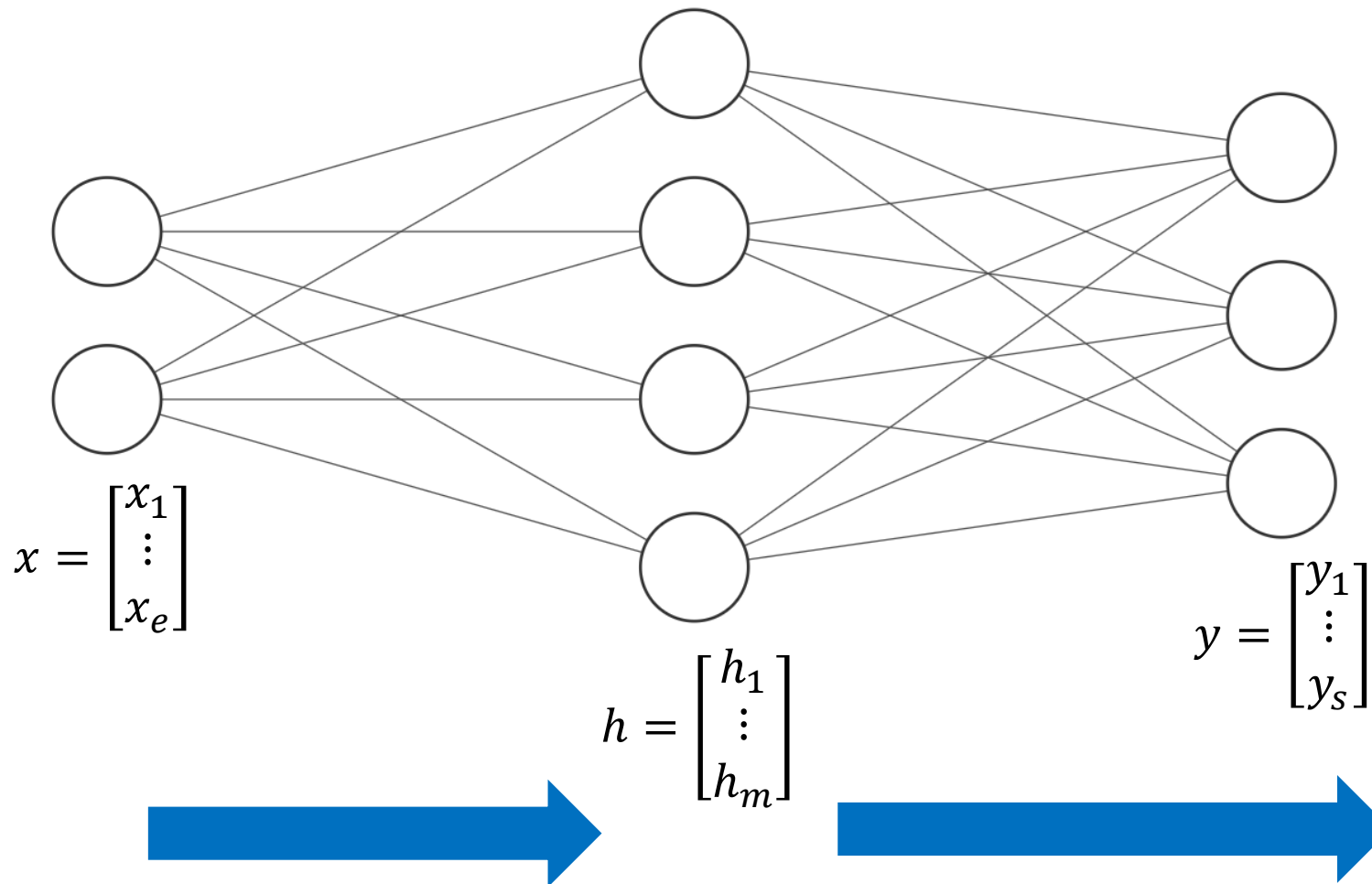
Veja a implementação do método no
Jupyter Notebook

Os objetivos são possibilitar uma implementação mais computacionalmente eficiente usando numpy e generalizar a solução para entradas e saídas com quaisquer tamanhos.

Implementando *Backpropagation* Matricialmente



Forward Pass

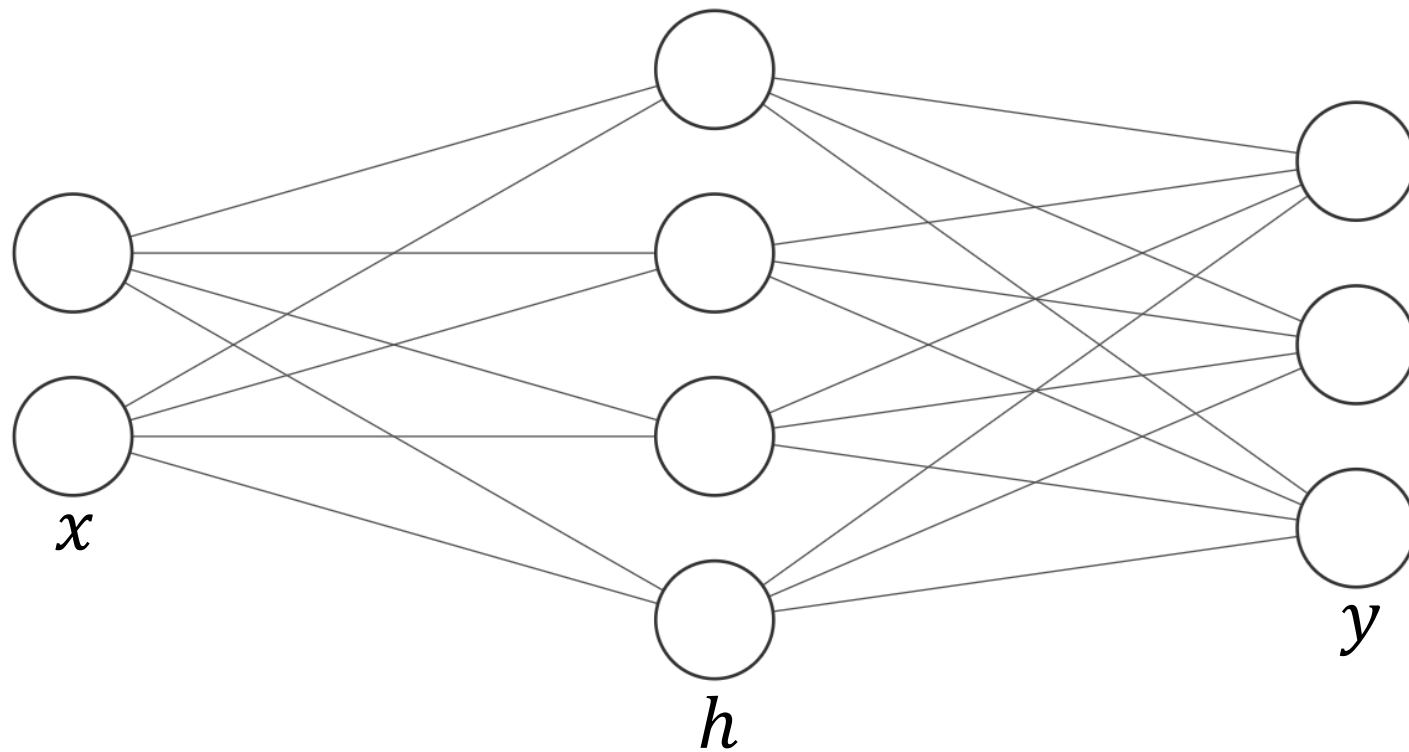


$$h = \sigma(W^h x + b^h)$$

$m \times 1$ $m \times e$ $e \times 1$ $m \times 1$

$$y = W^y h + b^y$$

$s \times 1$ $s \times m$ $m \times 1$ $s \times 1$



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (input), $m \times e$ (weights), $e \times 1$ (bias), $m \times 1$ (output).

$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (input), $s \times m$ (weights), $m \times 1$ (bias), $s \times 1$ (output).

$$\frac{\partial L}{\partial W^y} = \sum_{i=0}^N \left[\frac{2}{N} (y_i - \hat{y}_i) h^T \right]$$

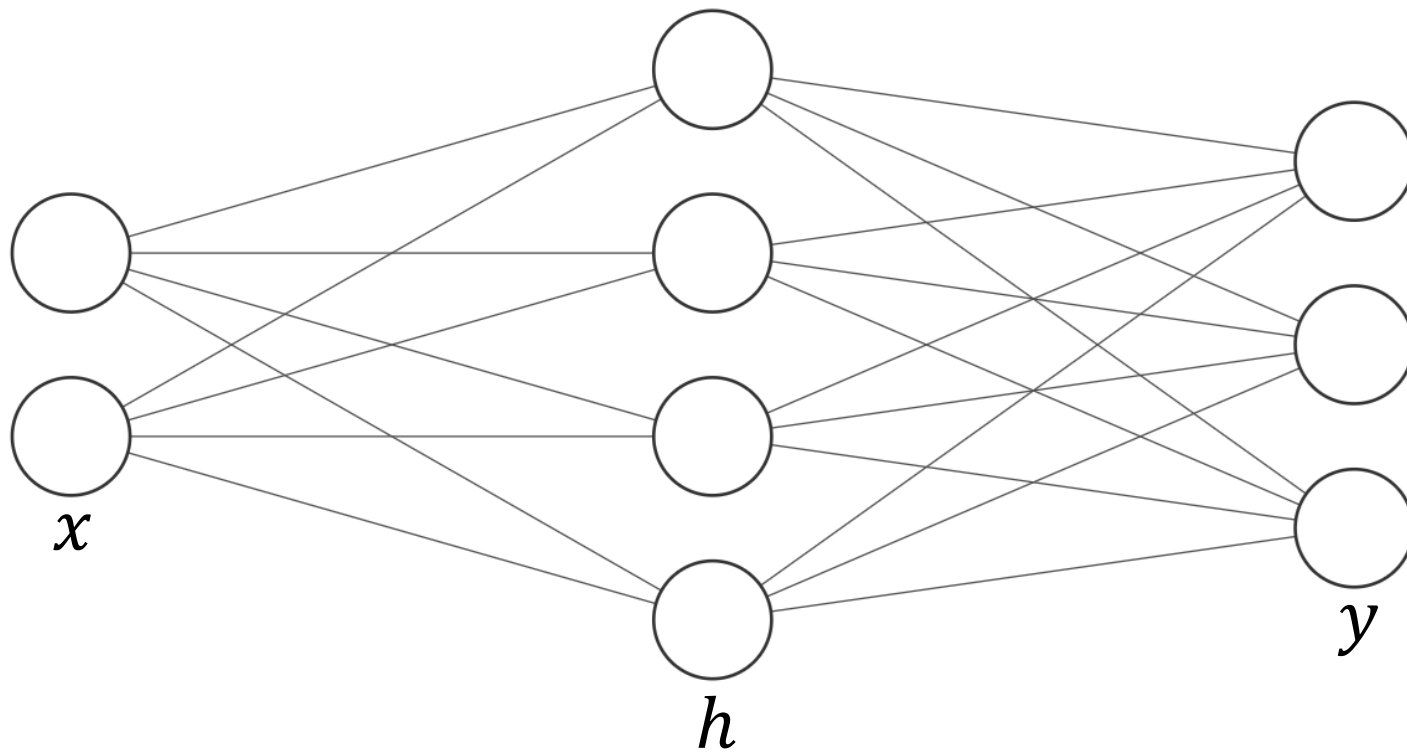
Dimensions: $s \times 1$ (error), $1 \times m$ (hidden state).

$\frac{\partial L}{\partial y} \quad \frac{\partial y}{\partial W^y}$

$$\frac{\partial L}{\partial b^y} = \sum_{i=0}^N \left[\frac{2}{N} (y_i - \hat{y}_i) \right]$$

Dimensions: $s \times 1$ (error).

Gradiente do erro em relação aos parâmetros da última camada



$$h = \sigma(W^h x + b^h)$$

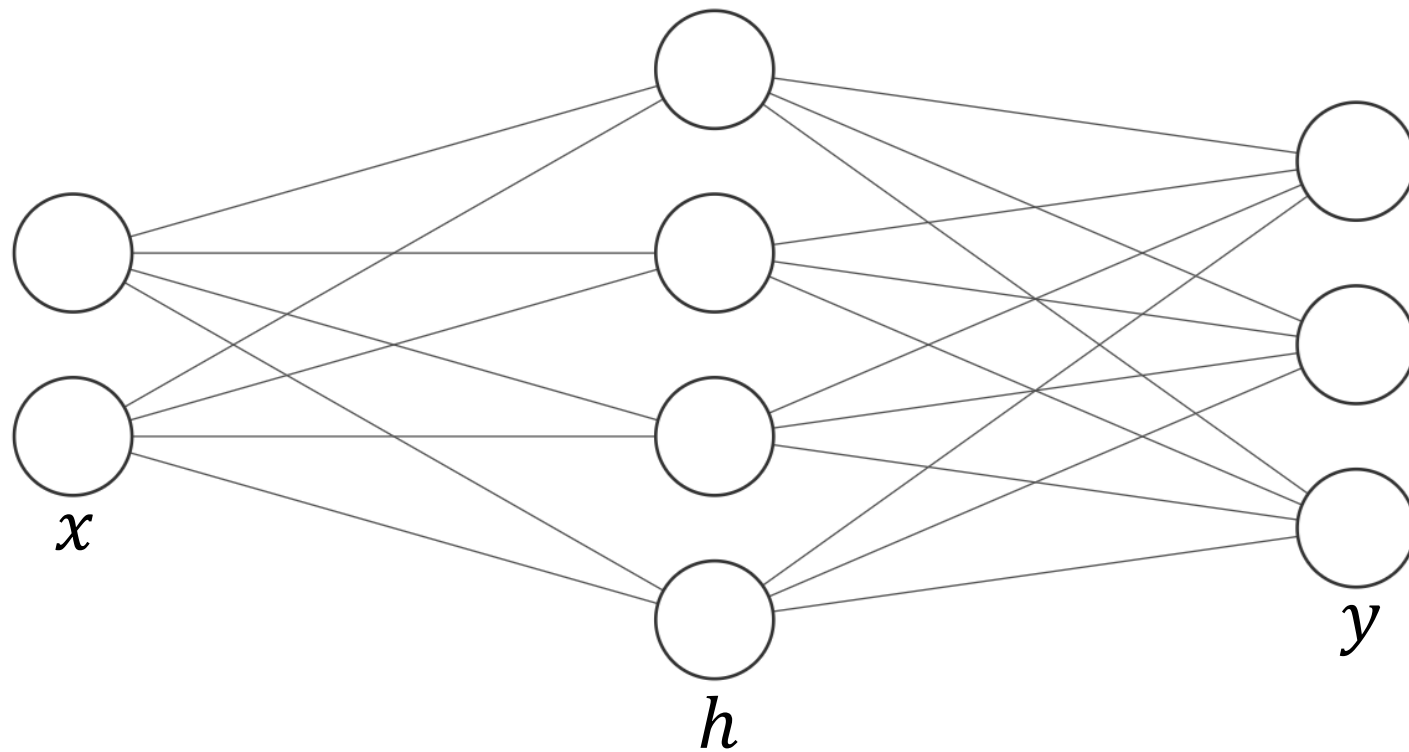
Dimensions: $m \times 1$ (for x), $m \times e$ (for W^h), $e \times 1$ (for b^h), and $m \times 1$ (for h).

$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (for y), $s \times m$ (for W^y), $m \times 1$ (for h), and $s \times 1$ (for b^y).

$$\frac{\partial L}{\partial W^h} = \sum_{i=0}^N \left[\sigma'(z^h) \odot \left(W^{y^T} \frac{2}{N} (y_i - \hat{y}_i) \right) \right] x^T$$

Gradiente do erro em relação aos parâmetros da camada oculta



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (for x), $m \times e$ (for W^h), $e \times 1$ (for b^h), and $m \times 1$ (for h).

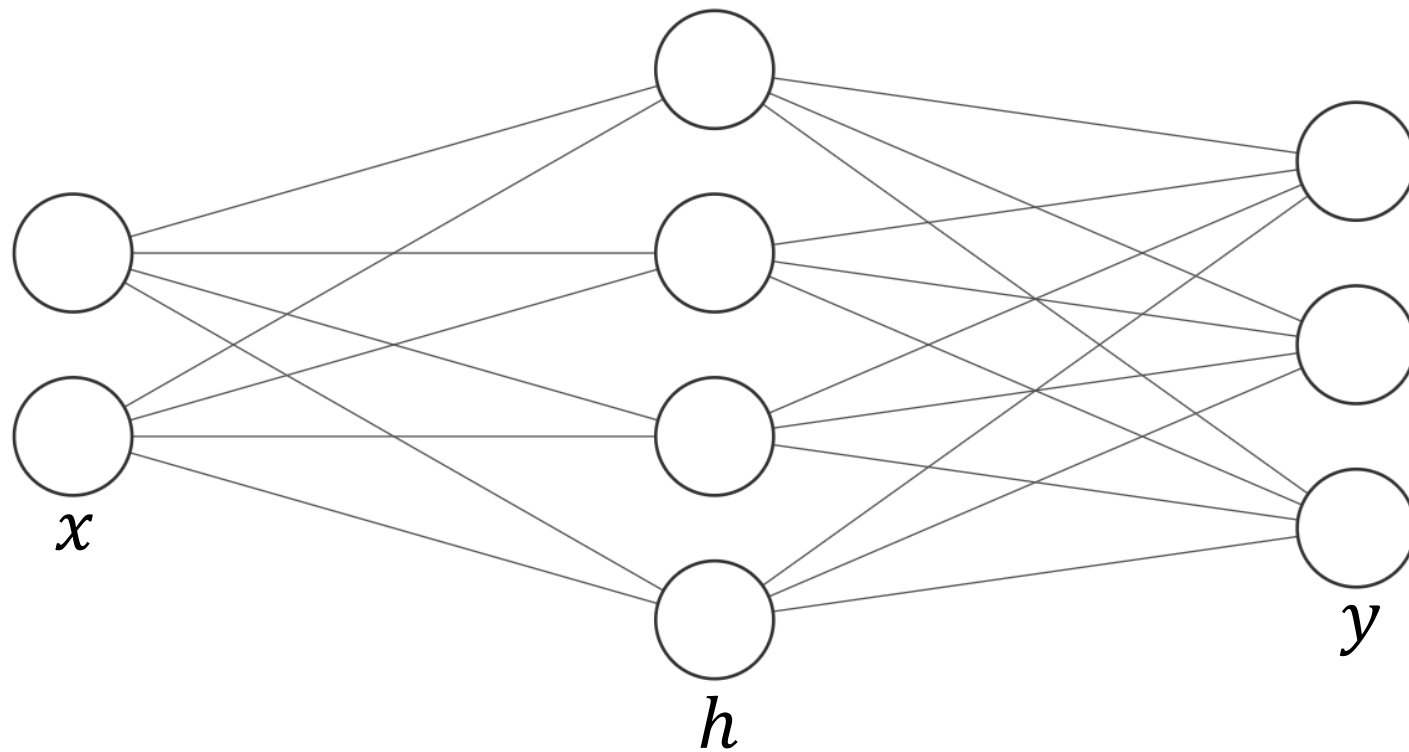
$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (for y), $s \times m$ (for W^y), $m \times 1$ (for h), and $s \times 1$ (for b^y).

$$\frac{\partial L}{\partial W^h} = \sum_{i=0}^N \left[\sigma'(z^h) \odot \underbrace{\left(W^{y^T} \frac{2}{N} (y_i - \hat{y}_i) \right)}_{\substack{m \times s \\ m \times 1}} \right] x^T$$

The term $\frac{\partial L}{\partial h}$ is indicated as having dimensions $m \times 1$.

Gradiente do erro em relação aos parâmetros da camada oculta



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (input), $m \times e$ (weights), $e \times 1$ (bias), $m \times 1$ (output).

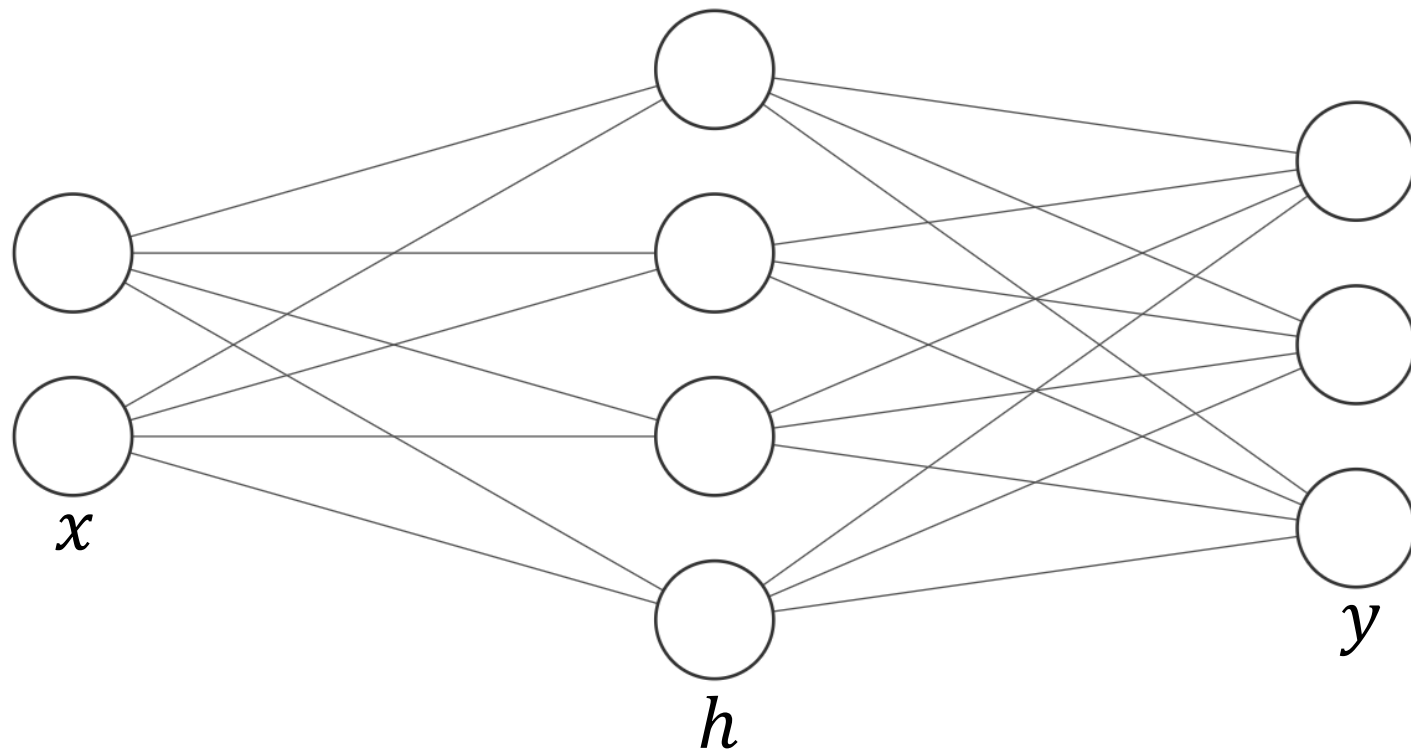
$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (input), $s \times m$ (weights), $m \times 1$ (bias), $s \times 1$ (output).

$$\frac{\partial L}{\partial W^h} = \sum_{i=0}^N \left[\underbrace{\sigma'(z^h) \odot \left(W^{yT} \frac{2}{N} (y_i - \hat{y}_i) \right)}_{\substack{\text{produto ponto a ponto} \\ m \times 1}} \right] x^T$$

Dimensions: $m \times s$ (weights), $s \times 1$ (output), $m \times 1$ (gradient).

Gradiente do erro em relação aos parâmetros da camada oculta



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (for x), $m \times e$ (for W^h), $e \times 1$ (for b^h), and $m \times 1$ (for h).

$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (for y), $s \times m$ (for W^y), $m \times 1$ (for h), and $s \times 1$ (for b^y).

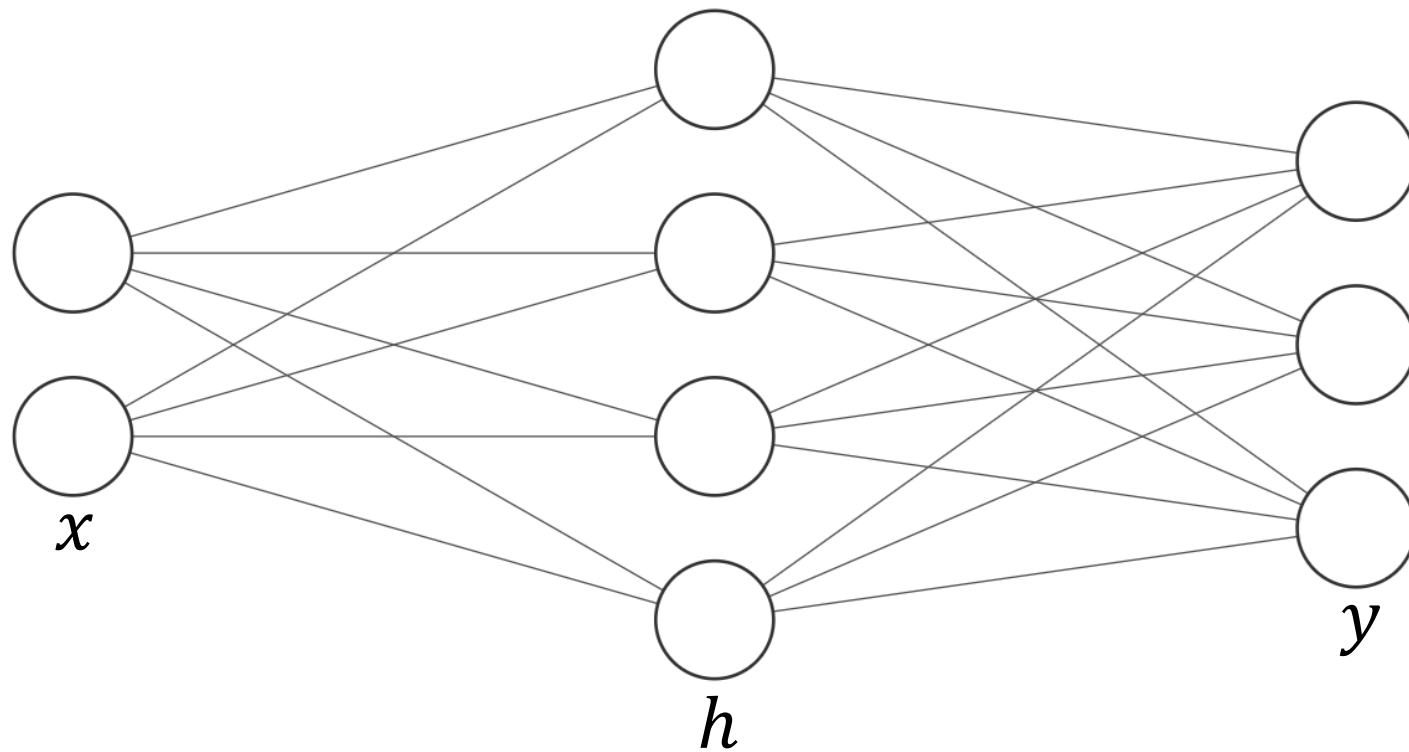
$$\frac{\partial L}{\partial W^h} = \sum_{i=0}^N \left[\underbrace{\sigma'(z^h) \odot \left(W^{yT} \frac{2}{N} (y_i - \hat{y}_i) \right)}_{m \times 1} \right] x^T$$

Dimensions: $m \times e$ (for $\frac{\partial L}{\partial W^h}$), $m \times s$ (for W^{yT}), $s \times 1$ (for $y_i - \hat{y}_i$), and $1 \times e$ (for x^T).

produto ponto a ponto (point-to-point product) is indicated for the operation \odot .

The term $\frac{\partial L}{\partial z^h}$ is indicated as the gradient of the error with respect to the hidden layer output, with dimension $m \times 1$.

Gradiente do erro em relação aos parâmetros da camada oculta



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (input), $m \times e$ (weights), $e \times 1$ (bias), $m \times 1$ (output).

$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (input), $s \times m$ (weights), $m \times 1$ (bias), $s \times 1$ (output).

analogamente...

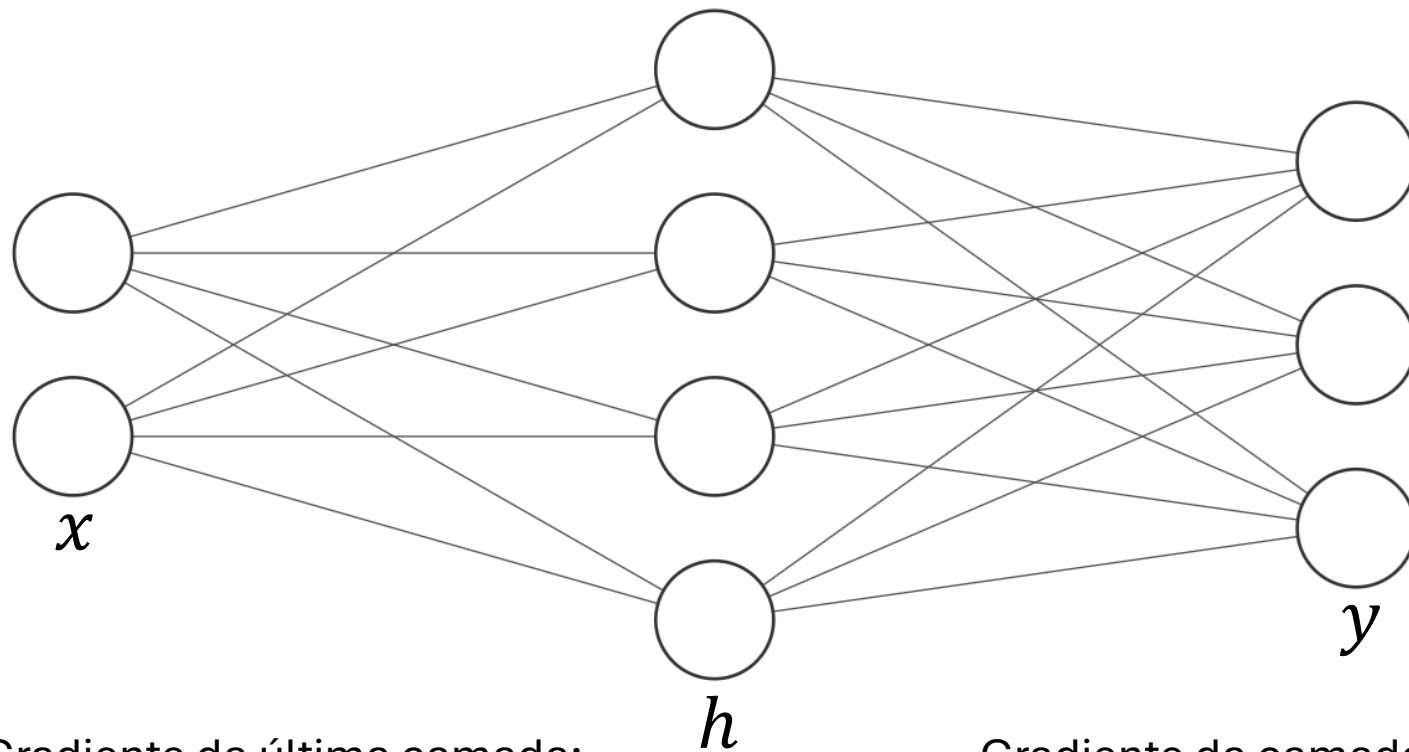
produto ponto a ponto

$$\frac{\partial L}{\partial b^h} = \sum_{i=0}^N \left[\underbrace{\sigma'(z^h) \odot \left(W^{yT} \frac{2}{N} (y_i - \hat{y}_i) \right)}_{m \times 1} \right]$$

Dimensions: $m \times 1$ (gradient), $m \times s$ (weights), $s \times 1$ (error).

Gradiente do erro em relação aos parâmetros da camada oculta

Resumo do Resultado Final



$$h = \sigma(W^h x + b^h)$$

Dimensions: $m \times 1$ (input x), $m \times e$ (weight matrix W^h), $e \times 1$ (bias vector b^h), and $m \times 1$ (hidden layer output h).

$$y = W^y h + b^y$$

Dimensions: $s \times 1$ (output y), $s \times m$ (weight matrix W^y), $m \times 1$ (hidden layer output h), and $s \times 1$ (bias vector b^y).

Gradiente da última camada:

$$\frac{\partial L}{\partial W^y} = \sum_{i=0}^N \left[\frac{2}{N} (y_i - \hat{y}_i) h^T \right]$$

$$\frac{\partial L}{\partial b^y} = \sum_{i=0}^N \left[\frac{2}{N} (y_i - \hat{y}_i) \right]$$

Gradiente da camada oculta:

$$\frac{\partial L}{\partial W^h} = \sum_{i=0}^N \left[\sigma'(z^h) \odot \left(W^{yT} \frac{2}{N} (y_i - \hat{y}_i) \right) \right] x^T$$

$$\frac{\partial L}{\partial b^h} = \sum_{i=0}^N \left[\sigma'(z^h) \odot \left(W^{yT} \frac{2}{N} (y_i - \hat{y}_i) \right) \right]$$

Veja a implementação do método no
Jupyter Notebook